# A Survey on Autonomous Driving Datasets: Statistics, Annotation Quality, and a Future Outlook

Mingyu Liu*⬤, Ekim Yurtsever⬤, *Member, IEEE*, Jonathan Fossaert⬤, Xingcheng Zhou⬤,
Walter Zimmer⬤, Yuning Cui⬤, *Student Member, IEEE*, Bare Luka Zagar⬤, Alois C. Knoll⬤, *Fellow, IEEE*

*Abstract*—**Autonomous driving has rapidly developed and shown promising performance due to recent advances in hardware and deep learning techniques. High-quality datasets are fundamental for developing reliable autonomous driving algorithms. Previous dataset surveys either focused on a limited number or lacked detailed investigation of dataset characteristics. To this end, we present an exhaustive study of 265 autonomous driving datasets from multiple perspectives, including sensor modalities, data size, tasks, and contextual conditions. We introduce a novel metric to evaluate the impact of datasets, which can also be a guide for creating new datasets. Besides, we analyze the annotation processes, existing labeling tools, and the annotation quality of datasets, showing the importance of establishing a standard annotation pipeline. On the other hand, we thoroughly analyze the impact of geographical and adversarial environmental conditions on the performance of autonomous driving systems. Moreover, we exhibit the data distribution of several vital datasets and discuss their pros and cons accordingly. Finally, we discuss the current challenges and the development trend of the future autonomous driving datasets.**
**https://github.com/MingyuLiu1/autonomous_driving_datasets.git**

*Index Terms*—**Dataset, impact score, data analysis, annotation quality, autonomous driving.**

## I. INTRODUCTION

**A**UTONOMOUS driving (AD) aims to revolutionize the transportation system by creating vehicles that can accurately perceive their environment, make intelligent decisions, and drive safely without human intervention. Due to thrilling technical development, various autonomous driving products have been implemented in several fields, such as robotaxis [1]. These rapid advancements in autonomous driving rely heavily on extensive datasets, which help autonomous driving systems be robust and reliable in complex driving environments.

In recent years, there has been a significant increase in the quality and variety of autonomous driving datasets. The first apparent phenomenon in the development of datasets is the various data collection strategies, including synthetic datasets [2]–[12] generated by simulators and recorded from

M. Liu, X. Zhou, W. Zimmer, Y. Cui, BL. Zagar, and AC. Knoll are with the Chair of Robotics, Artificial Intelligence and Real-Time Systems, Technical University of Munich, 85748 Garching bei München, Germany E-mail: {mingyu.liu, xingcheng.zhou, walter.zimmer, bare.luka.zagar}@tum.de, {yuning.cui, knoll}@in.tum.de

J. Fossaert is with the School of Engineering and Design, Technical University of Munich, 85748 Garching bei München, Germany (E-mail: jonathan.fossaert@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)
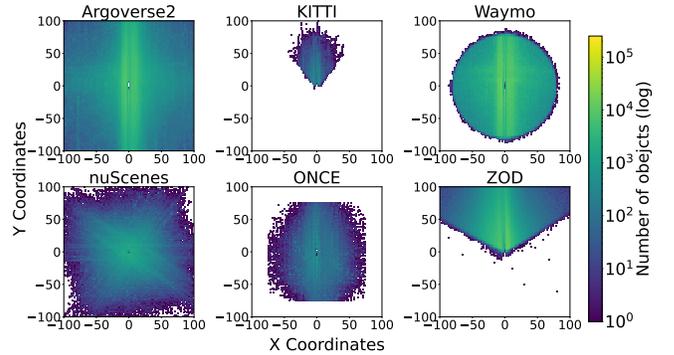
* Corresponding author

Fig. 1. Bird's-Eye View object distribution of datasets. Each heatmap represents a dataset and is plotted using X and Y coordinates. Y is the driving direction of the ego-vehicle. The unique annotation characters of each dataset are reflected in the distribution range, density, and number of bounding boxes.

the real world [13]–[29], to name just a few. Secondly, the datasets vary in composition, including but not limited to multiple sensor modalities like camera images and LiDAR point clouds, different annotation types for various tasks, and data distribution. Fig. 1 depicts the distribution of 3D object bounding box in six well-known real-world datasets (*Argoverse 2* [28], *KITTI* [13], *nuScenes* [22], *ONCE* [30], *Waymo* [23], and *ZOD* [31]) under a Bird's-Eye View (BEV), highlighting each dataset's unique annotation characteristics. The sensor mounting positions also reflect datasets' various sensing domains, including onboard, Vehicle-to-Everything (V2X), or drone domain. The datasets' geometric diversity and varying weather conditions also enhance the generalizability of autonomous driving datasets.

### A. Research Gap & Motivation

We demonstrate the yearly published number of perception datasets in Fig. 2 to illustrate the growth trends in autonomous driving datasets. Given the vast and growing number of publicly available datasets, a comprehensive survey of these resources is valuable for advancing both academic and industrial research in autonomous driving. In prior work, Yin et al. [32] summarized 27 publicly available datasets containing data collected on public roads. As a sequential work of [32], [33] extended the number of datasets. Guo et al. [34] and Janai et al. [35] proposed systematic introductions to the existing datasets from an application perspective. Beyond describing existing datasets, Liu et al. [36] discussed the domain adaptation between synthetic and real data and automatic labeling

TABLE I
WE COMPARE OUR SURVEY PAPER WITH OTHER AD DATASET SURVEYS IN THE FOLLOWING PERSPECTIVES: COLLECTED DATASET NUMBER (#DATASET), RELEVANT TASKS, SENSING DOMAIN (S. DOMAIN), SENSOR MODALITY (S. MODA.), GEOMETRIC CONDITIONS (GEO.), ENVIRONMENTAL CONDITIONS (ENV.), ANALYZING DATA DISTRIBUTION, INTRODUCING ANNOTATION QUALITY AND PROCESS. IN THE CONTEXT OF ENVIRONMENTAL CONDITIONS, WE REFER TO THE VARIABILITY IN WEATHER CONDITIONS AND ILLUMINATION. GEOMETRIC CONDITIONS INCLUDE SCENARIO TYPES AND GEOGRAPHICAL SCOPE. WE DESCRIBE THE TASK TYPES IN A COARSE GRANULARITY, INCLUDING PERCEPTION (PERC.), PREDICTION (PRED.), PLANNING (PL.), CONTROL (C.), 0AND END-TO-END (E2E).

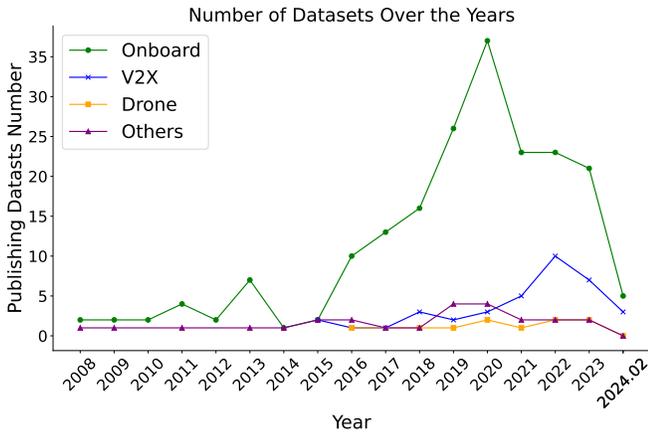| Survey | General | | | Data | | | | | Annotation | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | Year | #Datasets | Tasks | S. domain | S. moda. | Geo. | Env. | Data analysis | Quality | Process |
| When to use what dataset [32] | 2017 | 27 | Perc | | ✓ | ✓ | ✓ | | | |
| Self-driving Algorithm [33] | 2019 | 37 | Perc | | ✓ | ✓ | ✓ | | | |
| Is it safe to drive [34] | 2019 | 54 | Perc, Pred, E2E | | ✓ | ✓ | ✓ | | | |
| CV for AVs [35] | 2020 | 33 | Perc | | ✓ | ✓ | ✓ | | | ✓ |
| A Survey on AD Datasets [36] | 2021 | 30 | Perc | | ✓ | ✓ | ✓ | | | ✓ |
| 3D Semantic Segmentation [41] | 2021 | 29 | Perc | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| AD-Dataset [38] | 2022 | 204 | Perc, Pred, Pl, C | ✓ | ✓ | ✓ | | | | |
| Anomaly Detection [39] | 2023 | 16 | Perc | | ✓ | ✓ | ✓ | ✓ | | |
| Synthetic Datasets for AD [40] | 2023 | 17 | Perc, Pred | | ✓ | ✓ | ✓ | | | |
| Decision-making [42] | 2023 | 25 | Pl, C | | ✓ | ✓ | ✓ | | | |
| Open-sourced Data Ecosystem [37] | 2023 | 70 | Perc, Pred, Pl | ✓ | ✓ | ✓ | | | | ✓ |
| **Ours** | 2024 | 265 | Perc, Pred, Pl, C, E2E | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |



Fig. 2. Overview of dataset publication trends from 2008 to 2024. The diagram demonstrates a significant increase in the publication of onboard datasets between 2015 and 2020, followed by a gradual decline thereafter. In contrast, there has been a rising trend in the publication of V2X datasets, indicating growing research interest in cooperative perception systems.

methods. Li et al. [37] summarized existing datasets and undertook an exhaustive analysis of the characters of the next-generation datasets. However, these surveys only summarized a small number of datasets, causing a non-wide scope. AD-Dataset [38] collected a large number of datasets while lacking detailed analysis for the attributes of these datasets. Compared to studies on versatile datasets, some researchers presented surveys on a particular type of autonomous driving dataset, such as anomaly detection [39], synthetic datasets [40], 3D semantic segmentation [41], or decision-making [42]. Additionally, some task-specific surveys [43], [44] also organized the related AD datasets.

In this work, we present a comprehensive and systematic survey of a large number of datasets in autonomous driving. We compare our survey to others in Table I. Our survey covers all tasks from perception to control, considers real-world as well as synthetic data, and provides insights into the data modality and quality of several crucial datasets.

## B. Main Contributions

**The main contributions of this paper can be summarized as follows**:

- We present an overview of the most exhaustive survey on autonomous driving datasets recorded to date. We show publicly available datasets as comprehensively as possible, recording their fundamental characteristics, such as published year, data size, sensor modalities, sensing domains, geometrical and environmental conditions, and support tasks.
- We systematically illustrate the sensors and sensing domains for collecting AD data. Furthermore, we describe the main tasks in autonomous driving, including task goals, required data modalities, and evaluation metrics.
- We categorize datasets based on their sensing domains and relevant tasks, enabling researchers to efficiently identify and compile information for their target datasets. This approach facilitates more focused and productive research and development efforts.
- Additionally, we introduce an impact score metric to evaluate the influence of published perception datasets. This metric can also be a guide for developing future datasets. We deeply analyze datasets with high impact scores, highlighting their advantages and utility.
- We investigate the annotation quality of datasets and the existing labeling procedures for various autonomous driving tasks.
- Our detailed data statistics demonstrate the data distribution of various datasets from different perspectives, exhibiting their inherent limitations and suitable use cases.
- We analyze recent technology trends and the development direction of next-generation datasets, such as integrating language into AD data, generating AD data using Vision Language Models, standardizing data creation, and promoting an open data ecosystem.
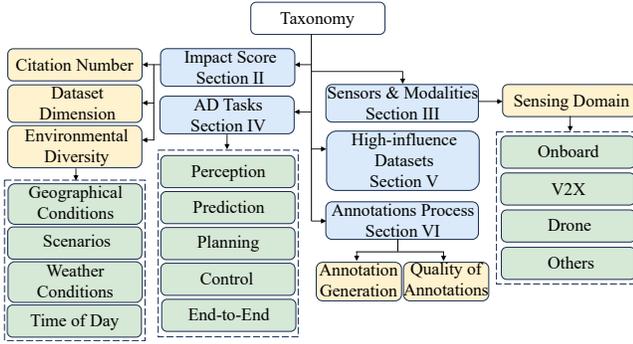
Fig. 3. This survey's primary taxonomy includes impact score, sensors and modalities, autonomous driving tasks, high-influence datasets, and annotation process.

### C. Scope & Limitations

We aim to conduct an exhaustive survey on the existing autonomous driving datasets to facilitate the development of future algorithms and datasets in this field. We collected datasets relevant to the five fundamental tasks of autonomous driving: perception, prediction, planning, control, and end-to-end (E2E) driving. To maintain clarity and prevent redundancy, we describe versatile datasets only within the main scope they support. Additionally, we collected a large number of datasets and exhibited them with their primary characters in tables. However, to ensure this survey aids researchers effectively, we focus our detailed discussions on the most impactful datasets rather than providing extensive descriptions of all datasets.

### D. Survey Structure

The rest of the survey is structured as follows: Section II introduces the approach leveraged to source public datasets and the evaluation metrics for them. Section III demonstrates the primary sensors used in autonomous driving and their modalities. Section IV discusses autonomous driving tasks, related challenges, and required data. In-depth discussions of several important datasets are presented in Section V. The process of annotations and factors affecting annotation quality are addressed in Section VI. Section VII provides statistical analysis of data distribution across various datasets. Future trends and potential research directions in autonomous driving datasets are explored in Section VIII. The paper concludes with Section IX. The survey's taxonomy is shown in Fig. 3.

## II. METHODOLOGY

This section comprises two parts: 1) the collection and filtering of datasets, as detailed in Section II-A, and 2) the evaluation metric of a dataset's impact on the autonomous driving domain, described in Section II-B.

### A. Datasets Collection

Following [45], we conducted a systematic review to exhaustively collect published autonomous driving datasets.

To ensure source diversity, we utilized well-known search engines such as *Google*, *Google Scholar* and *Baidu* to search datasets. To ensure a thorough dataset collection from various countries and regions, we conducted searches in English, Chinese, and German using keywords such as "autonomous driving dataset/benchmark", "intelligent vehicle dataset/benchmark" and terms related to object detection, classification, tracking, segmentation, prediction, planning, control, and end-to-end driving.

Furthermore, we explored *IEEE Xplore*, *Paperswithcode*, and pertinent conferences in autonomous driving and intelligent transportation systems to collate datasets from journals and conference proceedings. We verified datasets from these sources through keyword searches and manual title reviews.

Finally, to ensure the inclusion of specialized or lesser-known datasets, we searched through *Github* repositories. Similar to databases, we performed both manual and keyword-based searches for datasets.

### B. Dataset Evaluation Metrics

We introduce a novel metric, impact score, to assess the significance of a published dataset, which can also be a guide to preparing a new dataset. In this section, we explain in detail the approach to calculate the impact score of autonomous driving datasets.

For a fair and compatible comparison, we only consider datasets related to the perception domain, which takes up a large portion of autonomous driving datasets. Additionally, to ensure the objectivity and comprehensibility of our scoring system, we take into account various factors, including citation score, data dimension, and environmental diversity. All the values are gathered from official papers or open-source dataset websites.

**Citation Score.** First, we calculate the citation scores from the total citation number and average annual citation. To gain fair citation counts, we choose the time of the earliest version of a dataset as its publication time. Moreover, all citation counts were collected as of March 05, 2024, to ensure the comparison is based on a consistent timeframe. The total citation number $c^t$ reflects the overall influence of a dataset. A higher count of this metric means that the dataset has been widely recognized and utilized by researchers. However, datasets published in earlier years can accumulate more citations. To address this unfairness, we leverage average annual citation $c^a$, which describes a dataset's yearly citation increase speed. The function is shown in Eq. 1.

$$c^a = \begin{cases} c^t / (y_{curr} - y_{pub}) & \text{if } y_{curr} \neq y_{pub} \\ c^t & \text{if } y_{curr} = y_{pub} \end{cases} \quad (1)$$

where $y_{curr}$ and $y_{pub}$ represent the current year and the dataset published year, respectively. On the other hand, the citation number of distastes has a wide distribution range from single digits to tens of thousands. To alleviate the extreme imbalance and highlight the differences of each dataset, we apply a logarithmic transformation followed by Min-Max normalization to both $c^t$ and $c^a$, described in Eq. 2.

$$c_{norm} = \min\text{-}\max\left(\log(c)\right) \quad (2)$$

Finally, the citation score $c_{score}$ is the summation of $c_{norm}^t$ and $c_{norm}^a$:

$$c_{score} = 0.5c_{norm}^t + 0.5c_{norm}^a \qquad (3)$$

**Data Dimension Score.** We measure the data dimension across four perspectives: dataset size, temporal information, task number, and labeled categories. Dataset size $f$ is represented by the frame number of a dataset, reflecting its volume and comprehensiveness. To get the dataset size score $f_{norm}$, we leverage the same method as the citation score to process the frame number to overcome the extreme imbalance between different datasets.

Temporal information is essential for autonomous driving as it enables the vehicle to understand how the surrounding environment changes over time. We use $t \in \{0, 1\}$ to indicate whether a dataset includes temporal information. Regarding the task number, we only consider datasets related to the six fundamental tasks in the autonomous driving perception domain, such as 2D object detection, 3D object detection, 2D semantic segmentation, 3D semantic segmentation, tracking, and lane detection. Therefore, the task number score is recorded as $t^n \in \{1, 2, 3, 4, 5, 6\}$. A large number of categories is critical for the robustness and versatility of a dataset. During the statistic, if a dataset supports multiple tasks and includes various types of annotation, we choose the largest number of categories. Afterward, the categories are divided into five levels, $cat = \{1, 2, 3, 4, 5\}$, based on quintiles. We normalize $t_n$ and $l$ before the following process to simplify the calculation.

In order to reflect the data dimension score $d_{score}$ as objectively as possible, we give different weights to the four components, as shown in Eq. 4.

$$d_{score} = 0.5f_{norm} + 0.1t + 0.2t_{norm}^n + 0.2cat_{norm} \qquad (4)$$

**Environmental Diversity Score.** We evaluate the environmental diversity of a dataset according to the following factors: 1) weather conditions (such as rain or snow), 2) time of day (e.g., morning or dusk), 3) types of driving scenarios (like urban or rural), and 4) geometric scope refers to the number of countries or cities where the data is recorded. It is worth noting that we treat the geographical scope for synthetic datasets as undefined. We follow the granularity with which a paper categorizes its data to quantify the diversity. Moreover, for the missing value, if a dataset announces that the data is recorded under diversity conditions, we assign the median value; otherwise, the missing value is set to one. We apply quintiles to quantify each factor into five distinct levels. After that, the environmental diversity score $e_{score}$ is the sum of these four factors.

In the end, we leverage Eq. 5 to calculate the impact score $i_{score}$.

$$i_{score} = 60c_{score} + 20d_{score} + 20e_{score} \qquad (5)$$

The total impact score is 100, and 60 percent belongs to the citation score $c_{score}$. Data dimension score $d_{score}$ and environmental diversity score $e_{score}$ takes 40 percent.

## III. SENSORS AND PERCEPTION TECHNOLOGY IN AUTONOMOUS DRIVING

In this section, we introduce the sensors and their modalities mainly used in autonomous driving (III-A). Subsequently, in III-B, we analyze the data acquisition methods and cooperative perception technologies.
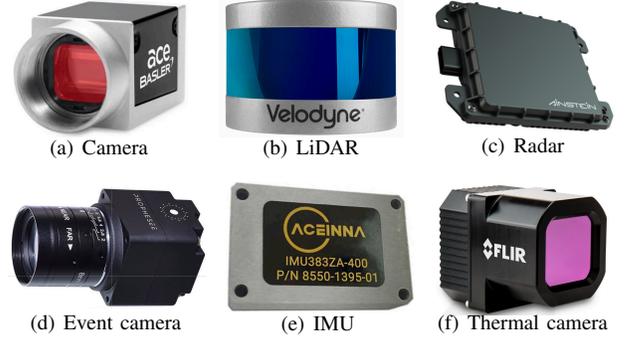
### A. Sensor Modalities



Fig. 4. Sensors on autonomous driving vehicles. The type of each sensor is (a) Camera: *Basler ace acA1600-20uc*, (b) LiDAR: *Velodyne Puck LITE*, (c) Radar: *Ainstein Launches K-79*, (d) Event-based camera: *Evaluation Kit 4 HD*, (e) IMU: *IMU383_Aceinna-W* and (f) Thermal camera: *FLIR_2nd_Gen_ADK*. All figures are extracted from the websites hosting the sensors.
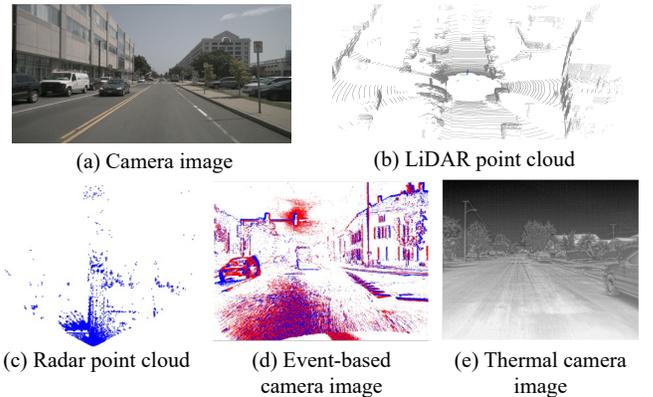


Fig. 5. We present the sensor modalities to provide an intuitive understanding of each sensor's characteristics. (a) is from nuScenes [22], (b) is from KITTI [13], (c) is from [46], (d) is from [47], (e) is from [48]. All figures are collected from the open-source data of datasets.

Efficiently and precisely collecting data from the surrounding environment is the key to a reliable perception system in autonomous driving. To achieve this goal, various types of sensors are utilized on self-driving vehicles and infrastructures. The sensor examples are shown in Fig. 4. The most used sensors are cameras, LiDARs, and radars. Event-based and thermal cameras are also installed on vehicles or roadside infrastructure to improve the perception capability further.

**RGB Images.** RGB images are usually recorded by monocular, stereo, or fisheye cameras. Monocular cameras offer a 2D view without depth; stereo cameras, with their dual lenses, provide depth perception; fisheye cameras use wide-angle lenses to capture a broad view. As shown in Fig. 5

(a), the 2D images capture color information, rich textures, patterns, and visual details of the environment. Due to these characters, RGB images are mainly used to detect vehicles and pedestrians and recognize road signs. However, the RGB images are vulnerable to conditions like low illumination, rain, fog, or glare [49].

**LiDAR Point Clouds.** LiDARs use laser beams to measure the distance between the sensor and an object, creating a 3D environment representation [50]. LiDAR point clouds (Fig. 5 (b)) provide precise spatial information with high resolution and can detect objects over long distances. However, the density of these points can decrease with increasing distance, leading to sparser representations for distant objects. Weather conditions, e.g., fog, also limit the performance of LiDAR. In general, LiDARs are suitable in cases that require 3D concise information.

**Radar Point Clouds.** Radars detect objects, distance, and relative speed by emitting radio waves and analyzing their reflection. Moreover, radars are robust under various weather conditions [51]. Nevertheless, radar point clouds are generally coarser than LiDAR data, lacking the detailed shape or texture information of objects. Therefore, radars are generally used to assist other sensors. Fig. 5 (c) exhibits the radar point clouds.

**Event of Event-based Camera.** Event-based cameras asynchronously capture data, activating only when a pixel detects a change in brightness. The captured data is called events (Fig. 5 (d)). Thanks to the specific data generation method, the recorded data has extremely high temporal resolution and captures fast motion without blur [52].

**Infrared Images of Thermal Camera.** Thermal cameras (see Fig. 5 (e)) detect heat signatures by capturing infrared radiation [53]. Due to producing images based on temperature differences, thermal cameras can work in total darkness and are unaffected by fog or smoke. However, thermal cameras cannot discern colors or detailed visual patterns evident. Furthermore, the resolution of infrared images is lower compared to optical cameras.

**Inertial Measurement Unit (IMU).** An IMU is an electronic device that measures and reports an object's specific force, angular rate, and sometimes magnetic field surrounding the object [54]. In autonomous driving, it is used to track the movement and orientation of the vehicle.

We analyze the sensor distribution from the collected datasets, shown in Fig. 6. More than half of the sensors are monocular cameras (52.79%) due to their low price and reliable performance. Additionally, 93 datasets (25.98%) include LiDAR data, valued for its high resolution and precise spatial information. However, its high-cost limits LiDAR's widespread use. Beyond LiDAR point clouds, 29 datasets leverage stereo cameras to capture depth information. Furthermore, 5.31%, 3.35%, and 1.68% datasets include radar, thermal camera, and fisheye camera. Given the temporal efficiency of capturing dynamic scenes, ten datasets generate data based on event-based cameras (2.79%).

### B. Sensing Domains and Cooperative Perception Systems

Sensory data and cooperation between the ego vehicle and surrounding entities are pivotal for ensuring autonomous
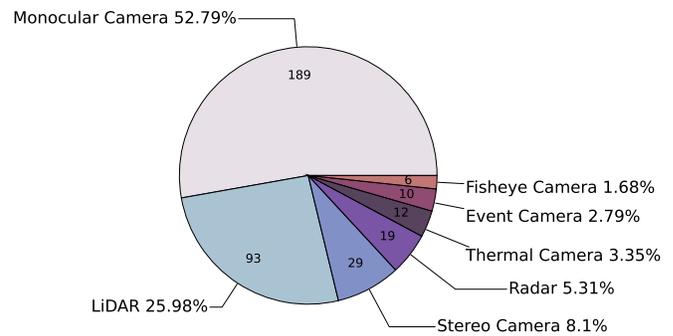


Fig. 6. Sensor type distribution. We show the distribution of different sensors. Overall, RGB cameras and LiDARs are the most used sensors in autonomous driving datasets.

driving systems' safety, efficiency, and overall functionality. Therefore, the positioning of sensors is crucial as it determines the quality, angle, and scope of data that can be collected. Generally, sensing domains in autonomous driving can be categorized as onboard, Vehicle-to-Everything (V2X), drone-based, and others.
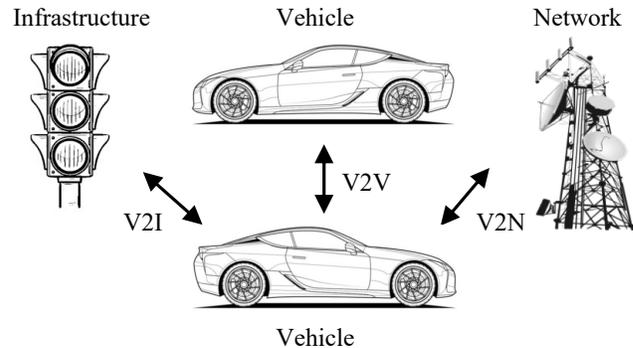


Fig. 7. Overview of cooperative perception systems in autonomous driving. A complete autonomous driving perception system consists of the ego vehicle along with collaborative interactions between vehicles, infrastructure, and networks.

**Onboard.** Onboard sensors are installed directly on the autonomous vehicle and usually consist of cameras, LiDARs, radars, and IMUs. These sensors provide a direct perspective from the vehicle's standpoint, offering immediate feedback on the surroundings. Nevertheless, due to the limited detection scope, onboard sensors may have limitations in providing advanced warnings about obstacles in blind spots or detecting hazards around sharp bends.

**V2X.** Vehicle-to-Everything encompasses communications between a vehicle and any other components in the transport system [55], including Vehicle-to-Vehicle, Vehicle-to-Infrastructure, and Vehicle-to-Network (as shown in Fig. 7). Beyond the immediate sensory input, the cooperative systems ensure multiple entities work harmoniously.

- *Vehicle-to-Vehicle (V2V)*
  V2V enables nearby vehicles to share essential data, such as position, velocity, and captured sensory data (e.g., camera images or LiDAR scans). This shared information

contributes to a more comprehensive understanding of driving scenarios.

- *Vehicle-to-Infrastructure (V2I)*
  V2I facilitates communications between the autonomous vehicle and infrastructure components such as traffic lights, signs, or roadside sensors. The sensors implemented in road infrastructure collaborate to enhance the perception range and situational awareness of autonomous vehicles. In this survey, we categorize interactions involving single or multiple vehicles with single or multiple infrastructure components, as well as collaborations among multiple infrastructure elements, under V2I.

- *Vehicle-to-Network (V2N)*
  V2N refers to exchanging information between a vehicle and a broader network infrastructure, often leveraging cellular networks to provide vehicles with access to cloud data. V2N aids the cooperation perception of V2V and V2I by sharing cross-area data or offering real-time updates about traffic congestion or road closures.

**Drone.** Drones, or Unmanned Aerial Vehicles (UAVs), offer an aerial perspective, providing data essential for trajectory prediction and route planning [18]. For example, the real-time data from drones can be integrated into traffic management systems to optimize the traffic flow and alert autonomous vehicles of accidents ahead.

**Others.** Data not collected by the previous three types is defined as others, such as other devices installed on non-vehicle objects or multiple domains.

## IV. Tasks in Autonomous Driving

This section offers insight into the pivotal tasks in autonomous driving, such as perception and localization (IV-A), prediction (IV-B), and planning and control (IV-C). The overview of the autonomous driving pipeline is demonstrated in Fig. 8. We detail their objectives, the nature of data they rely upon, and inherent challenges. Fig. 9 highlights several main tasks in autonomous driving.

### A. Perception and Localization

Perception focuses on understanding the environment based on the sensory data, while localization determines the autonomous vehicle's position within that environment.

**2D/3D Object Detection.** 2D or 3D object detection aims to identify the locations and classification of other entities within the driving environment. Although the detection technologies have significantly advanced, several challenges remain, such as object occlusions, varying light conditions, and diverse object appearances.

Usually, the Average Precision (AP) metric [61] is applied to evaluate the object detection performance. According to [1], the AP metric can be formulated as

$$AP = \int_0^1 \max \{ p(r'|r' \geq r) \} \, dr \qquad (6)$$

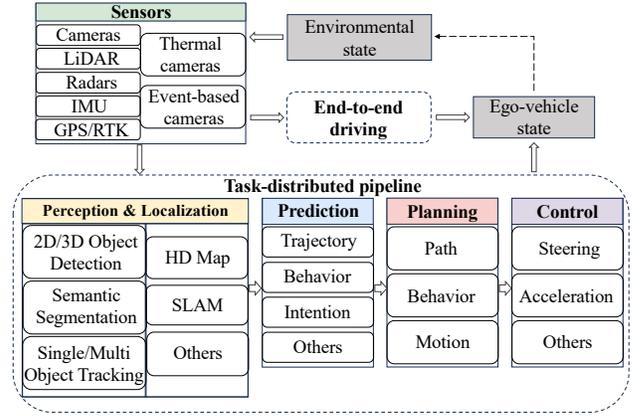where $p(r)$ is the precision-recall curve.



Fig. 8. The overview of autonomous driving pipeline. Autonomous driving systems can be categorized into two types: modular-based and end-to-end. Both types rely on data collected by various sensors installed on vehicles or infrastructures. These systems interact with and respond to the surrounding environment during driving scenarios.

**2D/3D Semantic Segmentation.** Semantic segmentation involves classifying each pixel of an image or point of a point cloud to its semantic category. From a dataset perspective, maintaining fine-grained object boundaries while managing extensive labeling requirements presents significant challenges for this task.

As mentioned in [62], the main metrics used for segmentation are mean Pixel Accuracy (mPA):

$$mPA = \frac{1}{k+1} \sum_{i=0}^{k} \frac{p_{ii}}{\sum_{j=0}^{k} p_{ij}} \qquad (7)$$

And the mean Intersection over Union (mIoU):

$$mIoU = \frac{1}{k+1} \sum_{i=0}^{k} \frac{p_{ii}}{\sum_{j=0}^{k} p_{ij} + \sum_{j=0}^{k} p_{ji} - p_{ii}} \qquad (8)$$

Where $k \in \mathbb{N}$ is the number of classes, and $p_{ii}$, $p_{ij}$, and $p_{ji}$ represent true positives, false positives, and false negatives, respectively.

**Object Tracking.** Object Tracking monitors the trajectories of a single or multiple objects over time. This task necessitates time-series RGB data, LiDAR, or radar sequences. Usually, object tracking includes single-object tracking or multi-object tracking (MOT).

Multi-Object-Tracking Accuracy (MOTA) is a widely utilized metric for multiple object tracking, which combines false negatives, false positives, and mismatch rate [63] (see Eq. 9).

$$MOTA = 1 - \frac{\sum_t (fp_t + fn_t + e_t)}{\sum_t gt_t} \qquad (9)$$

where $fp$, $fn$, and $e$ are the number of false positives, false negatives, and mismatch errors over time $t$. $gt$ is the ground truth.

Furthermore, instead of considering a single threshold, Average MOTA (AMOTA) is calculated based on all object confidence thresholds [64].

**HD Map.** HD mapping aims to construct detailed, highly accurate representations that include information about road
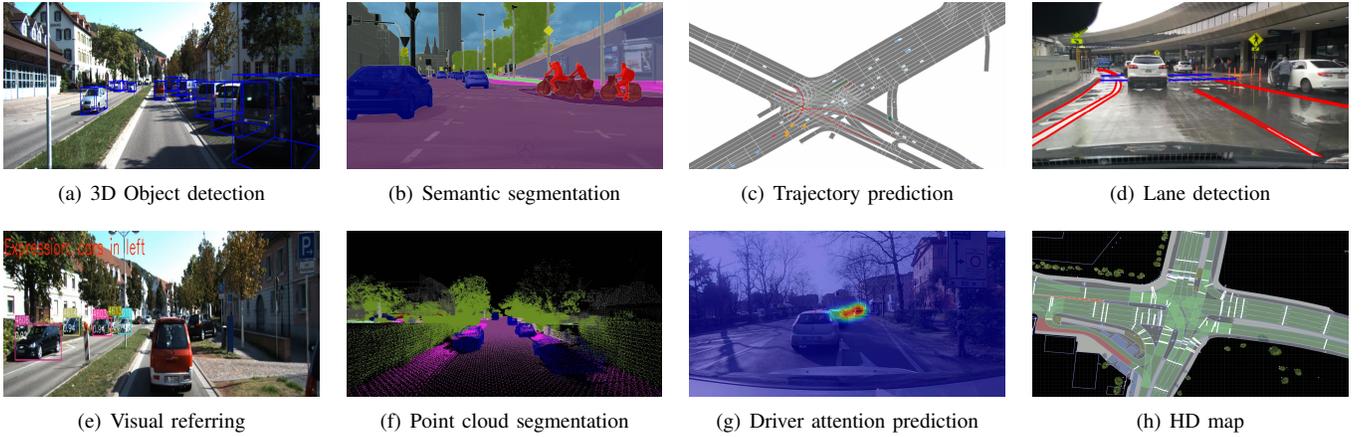
Fig. 9. Examples of various autonomous driving tasks. (a) is from KITTI [13], (b) is from Cityscapes [14], (c) is from V2X-Seq [56], (d) is from BDD100K [24], (e) is from Refer-KITTI [57], (f) is from KITTI-360 [58], (g) is from Dr(eye)ve [59], (h) is from TUMTraf [60]. All figures are collected from the open-source data of datasets or the websites hosting the datasets.

structures, traffic signs, and landmarks. A dataset should provide LiDAR data for precise spatial information and camera data for visual details to ensure established map accuracy. According to [28], HD map automation [65] and HD map change detection [66] have received more and more attention. Usually, the HD map quality is estimated using the accuracy metric.

**SLAM.** Simultaneous Localization And Mapping (SLAM) entails building a concurrent map of the surrounding environment and localizing the vehicle within this map. Hence, data from cameras, IMUs for position tracking, and real-time LiDAR point clouds are vital. Sturm et al. [67] introduces two evaluation metrics, relative pose error (RPE) and absolute trajectory error (ATE), for evaluating the quality of the estimated trajectory from the input RGB-D images.

### B. Prediction

Prediction refers to forecasting the future states or actions of surrounding agents. This capacity ensures safer navigation in dynamic environments. Several evaluation metrics are used for prediction [68], [69], such as Root Mean Squared Error (RMSE):

$$RMSE = \sqrt{\frac{1}{N}\sum_{n=1}^{N}\left(T_{pred}^n - T_{gt}^n\right)^2} \quad (10)$$

where $N$ is the total number of samples, $T_{pred}$ and $T_{gt}$ represent the predicted trajectory and ground truth.

Negative Log Likelihood (NLL) (see Eq. 11) is another metric focusing on determining the correctness of the trajectory, which can be used to compare the uncertainty of different models [70].

$$NLL = -\sum_{c=1}^{C} n_c log\left(\hat{n}_c\right) \quad (11)$$

Where $C$ is the total classes, $n_c$ is the binary indicator of the correctness of prediction, and $\hat{n}_c$ is the corresponding prediction probability.

**Trajectory Prediction.** Based on time-series data from sensors like cameras and LiDARs, trajectory prediction pertains to anticipating the future paths or movement patterns of other entities [68], such as pedestrians, cyclists, or other vehicles.

**Behavior Prediction.** Behavior prediction anticipates the potential actions of other road users [69], e.g., whether a vehicle will change the lane. Training behavior prediction models rely on extensive annotated data due to entities' vast range of potential actions within various scenarios.

**Intention Prediction.** Intention prediction focuses on inferring the higher-level goals behind the actions of objects, involving a deeper semantic comprehension of the physical or mental activities of humans [71]. Because of the task's complexity, it requires data from perception sensors like cameras, traffic signals, and hand gestures.

### C. Planning and Control

*1) Planning:* Planning represents the decision-making process in reaction to the perceived environment and predictions. A classic three-level hierarchical planning framework comprises path, behavioral, and motion planning [59].

**Path Planning** Path planning, also known as route planning, involves setting long-term objectives. It is a high-level process of determining the best path to the destination.

**Behavior Planning.** Behavior planning sits at the mid-level of the framework and is related to decision-making, including lane changes, overtaking, merging, and intersection crossing. This process relies on the correct understanding and interaction with the behavior of other agents.

**Motion Planning.** Motion planning deals with the actual trajectory the vehicle should follow in real time, considering obstacles, road conditions, and the predicted behavior of other road agents. In contrast to path planning, motion planning generates appropriate paths to achieve local objectives [59].

*2) Control:* Control mechanisms in autonomous driving govern how the self-driving car executes the decided path or behavior from the motion planning system and corrects tracking errors [72]. It translates high-level commands into actionable throttle, brake, and steering commands.
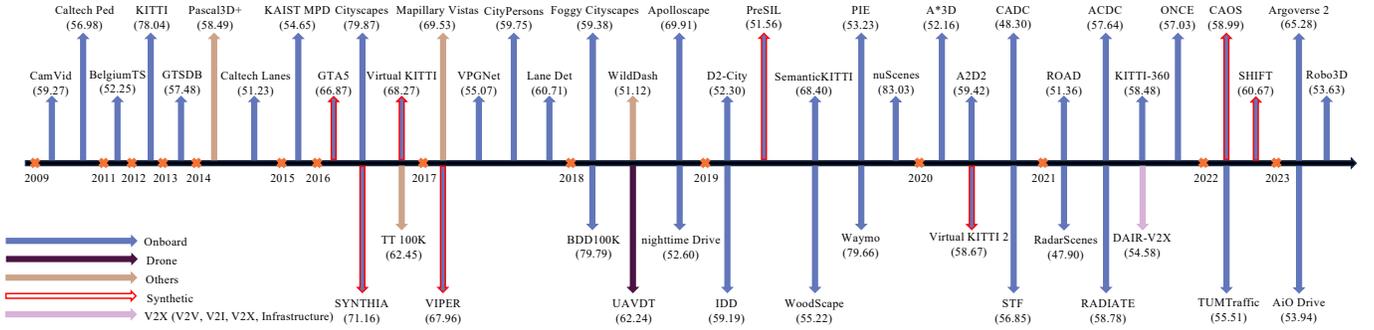
**Fig. 10.** Chronological overview. We show the top 50 datasets ranked by their impact scores. This diagram covers datasets released from 2009 through March 17, 2024.

## D. End-to-end Autonomous Driving

End-to-end approaches in autonomous driving are those where a single deep learning model handles everything from perception to control, bypassing the traditional modular pipeline. Such models can often be more adaptive as they rely on adjusting the whole model through learning. Their inherent promise is in their simplicity and efficiency by reducing the need for hand-crafted components [73]. However, implementing end-to-end models faces limitations such as extensive training data requirements, low interpretability, and inflexible modular tuning.

Large-scale benchmarking for end-to-end AD can be categorized into closed-loop and open-loop evaluation [74]. The closed-loop evaluation is based on a simulated environment, while open-loop evaluation involves assessing a system's performance against expert driving behavior from real-world datasets.

## V. High-Influence Datasets

This section describes the milestone autonomous driving datasets in the field of perception (V-A), prediction, planning, and control (V-B). We also exhibit datasets for end-to-end autonomous driving (V-C).

### A. Perception Datasets

Perception datasets are critical for developing and optimizing autonomous driving systems. They enhance vehicle reliability and robustness by providing rich, multimodal sensory data, ensuring effective perception and understanding of surroundings.

We leverage the proposed dataset evaluation metrics (see II-B) to calculate the impact scores of the collected perception datasets, subsequently selecting the top 50 datasets based on these scores to create a chronological overview, shown in Fig. 10. Concurrently, as described in III-B, we categorize datasets into onboard, V2X, drone, and others, choosing a subset from each to compile a comprehensive table of 50 datasets (Tab. II). It is worth noting that the datasets in the table are sorted by impact score within their respective categories and do not represent the overall top 50. In the following section, we choose several datasets with the highest impact scores within each sensing source and exhibit them, considering their published year.

*1) Onboard:* **KITTI.** KITTI [13] has been instrumental in advancing the autonomous driving domain since its release in 2012. KITTI data is recorded by various sensors, including cameras, LiDAR, and GPS/IMU, facilitating the development of algorithms for object detection, tracking, optical flow, depth estimation, and visual odometry. However, KITTI data was mainly recorded under ideal weather conditions in German urban areas. The relatively small geographic and environmental conditions scope limit its real-world applications.

*Cityscapes.* Cityscapes [14] comprises a vast collection of images captured explicitly in intricate urban environments and has become a standard benchmark semantic segmentation. Cityscapes provides pixel-level segmentation through meticulous labeling for 30 object classes, including various vehicle types, pedestrians, roads, and traffic sign information. However, Cityscapes primarily focuses on German cities and lacks diverse weather conditions. These limitations hinder its usefulness for developing robust and generalized autonomous driving solutions.

*VIPER.* VIPER [5] is a synthetic dataset collected based on driving, riding, and walking perspectives in a realistic virtual world. VIPER contains over 250K video frames annotated with ground-truth data for low- and high-level vision tasks. It encapsulates various weather conditions, lighting scenarios, and complex urban landscapes, making it an invaluable resource for testing the robustness of autonomous driving algorithms. Nonetheless, the transition from synthetic data to real-world application remains challenging, as algorithms must bridge the domain gap to maintain their effectiveness in real environments.

*SemanticKITTI.* SemanticKITTI [6] consists of over 43,000 LiDAR point cloud frames, making it one of the most extensive datasets for 3D semantic segmentation in outdoor environments. SemanticKITTI provides precise labels for 28 categories, such as car, road, building, etc., achieving a robust benchmark for semantic segmentation. However, similar to the previous datasets, SemanticKITTI faces limitations in real-world applicability due to its restricted environmental diversity and geographical scope.

*nuScenes.* nuScenes [22] stands as an essential contribution to the field of autonomous driving, providing multimodal sensor setup, including LiDAR, radars, and cameras. It addresses the diversity in urban scenes and environmental conditions.

TABLE II

HIGH-IMPACT PERCEPTION DATASETS. FOR A MORE COMPREHENSIVE DEMONSTRATION, WE EXHIBIT 50 PERCEPTION DATASETS FROM DIFFERENT SENSING DOMAINS INSTEAD OF THOSE WITH THE HIGHEST SCORES. THE CITATION NUMBER OF EACH DATASET IS NOT SHOWN IN THE TABLE DUE TO THE FAST CHANGE FOR THIS POINT.

| Dataset | Year | Size | Temp | Tasks — 2D Det | 3D Det | 2D Seg | 3D Seg | Tracking | Lane Det | Categories number | Weather conditions | Time of day | Scenario type | Geographical scope | Impact score |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Onboard** | | | | | | | | | | | | | | | |
| nuScenes [22] | 2019 | 40K | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | 23 | 3 | 2 | 4 | 2 | 83.03 |
| Cityscapes [14] | 2016 | 25K | × | ✓ | | ✓ | | | | 30 | 1 | 1 | 3 | 1 | 79.87 |
| BDD100K [24] | 2020 | 12M | ✓ | ✓ | | ✓ | | ✓ | ✓ | 40 | 5 | 2 | 4 | 1 | 79.79 |
| Waymo [23] | 2019 | 230K | ✓ | ✓ | ✓ | ✓ | | ✓ | | 23 | 2 | 3 | 5 | 1 | 79.66 |
| KITTI [13] | 2012 | 41K | ✓ | ✓ | ✓ | ✓ | | | | 8 | 1 | 1 | 2 | 1 | 78.04 |
| SYNTHIA [2] | 2016 | 13.4K | × | | | ✓ | | | | 13 | 2 | 3 | 5 | - | 71.16 |
| Apolloscape [17] | 2018 | 143,906 | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | 28 | 4 | 3 | 3 | 1 | 69.91 |
| SemanticKITTI [6] | 2019 | 43,552 | ✓ | | | | ✓ | | | 28 | 1 | 1 | 4 | 1 | 68.40 |
| Virtual KITTI [3] | 2016 | 21,260 | ✓ | | ✓ | | | ✓ | | 8 | 3 | 2 | 5 | - | 68.27 |
| VIPER [5] | 2017 | 254,064 | ✓ | ✓ | ✓ | ✓ | | ✓ | | 11 | 4 | 3 | 4 | - | 67.96 |
| GTA5 [4] | 2016 | 24,966 | × | | | ✓ | | | | 19 | 2 | 2 | 2 | - | 66.87 |
| Argoverse 2 [28] | 2023 | 6M | ✓ | ✓ | | ✓ | | ✓ | | 30 | 2 | 1 | 1 | 6 | 65.28 |
| Lane Det [75] | 2017 | 133,235 | × | | | | | | ✓ | 1 | 1 | 3 | 3 | 1 | 60.71 |
| SHIFT [10] | 2022 | 2,5M | ✓ | ✓ | ✓ | ✓ | | ✓ | | 23 | 5 | 5 | 3 | - | 60.67 |
| CityPersons [76] | 2017 | 25K | × | ✓ | | ✓ | | | | 30 | 2 | 1 | - | 27 | 59.75 |
| A2D2 [25] | 2020 | 41,277 | × | ✓ | ✓ | ✓ | | | | 38 | 2 | 1 | 3 | 3 | 59.42 |
| Foggy Cityscapes [77] | 2018 | 20,550 | × | ✓ | | ✓ | | | | 19 | 1 | 1 | 2 | 1 | 59.38 |
| CamVid [78] | 2009 | 701 | × | | | ✓ | | | | 32 | 2 | 1 | 2 | 1 | 59.27 |
| IDD [79] | 2019 | 10,004 | × | | | ✓ | | | | 34 | 3 | 3 | 5 | 1 | 59.18 |
| CAOS [7] | 2022 | 13K | × | | | ✓ | | | | 13 | 3 | 2 | 3 | - | 58.98 |
| RADIATE [26] | 2021 | 44,140 | ✓ | | ✓ | | | ✓ | | 8 | 5 | 2 | 4 | 1 | 58.78 |
| Virtual KITTI 2 [9] | 2020 | 20,992 | ✓ | ✓ | ✓ | ✓ | | ✓ | | 8 | 4 | 2 | 3 | 1 | 58.67 |
| KITTI-360 [58] | 2021 | 150K | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | 37 | 1 | 1 | 1 | 1 | 58.48 |
| Dr(eye)ve [59] | 2018 | 555000 | × | | | ✓ | | | | 10 | 3 | 3 | 3 | - | 58.13 |
| ACDC [80] | 2021 | 4,006 | × | | | ✓ | | | | 19 | 4 | 2 | 3 | 1 | 57.64 |
| GTSDB [81] | 2013 | 900 | × | | ✓ | | | | | 4 | 2 | 2 | 3 | 1 | 57.48 |
| ONCE [30] | 2021 | 1M | × | ✓ | | ✓ | | | | 5 | 3 | 2 | 5 | 1 | 57.03 |
| Caltech Ped [82] | 2009 | 250K | × | ✓ | | | | | | 1 | 1 | 1 | 1 | 2 | 56.98 |
| STF [83] | 2020 | 13,500 | × | | | ✓ | | | | 1 | 4 | 2 | 3 | 4 | 56.85 |
| **V2X** | | | | | | | | | | | | | | | |
| TUMTraf [60], [84]–[86] | 2022 | 50,253 | ✓ | ✓ | ✓ | ✓ | | ✓ | ✓ | 10 | 6 | 5 | 10 | 1 | 55.51 |
| DAIR-V2X [27] | 2021 | 71,254 | ✓ | ✓ | ✓ | ✓ | | | | 10 | 2 | 2 | 2 | 1 | 54.58 |
| V2XSet [11] | 2022 | 11,447 | ✓ | ✓ | ✓ | ✓ | | | | 1 | 1 | 1 | 1 | - | 49.84 |
| V2V4Real [87] | 2023 | 40K | ✓ | ✓ | ✓ | | | ✓ | | 5 | 2 | 1 | 2 | 1 | 40.28 |
| Rope3D [88] | 2022 | 50K | × | ✓ | ✓ | ✓ | | | | 12 | 3 | 3 | 2 | 1 | 48.96 |
| V2X-Sim [89] | 2022 | 10K | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | 23 | 1 | 1 | 1 | - | 48.50 |
| V2VNet [90] | 2020 | 51,2K | ✓ | | ✓ | | | | | 1 | 1 | 1 | 1 | - | 48.20 |
| T&J [91] | 2019 | 100 | ✓ | ✓ | ✓ | ✓ | | | | 1 | 1 | 1 | 2 | 1 | 46.63 |
| Co-Percep [92] | 2020 | 10K | × | ✓ | | ✓ | | | | 1 | 1 | 1 | 1 | 1 | 42.99 |
| DeepAccident [12] | 2023 | 285K | ✓ | | | | | | | 6 | 4 | 3 | 2 | - | 41.79 |
| LUMPI [93] | 2022 | 200K | ✓ | ✓ | | ✓ | | ✓ | | 3 | 6 | 3 | 1 | 1 | 40.42 |
| **Drone** | | | | | | | | | | | | | | | |
| UAVDT [94] | 2018 | 80K | ✓ | ✓ | | | | ✓ | | 3 | 2 | 2 | 6 | 1 | 61.63 |
| DroneVehicle [95] | 2021 | 28,439 | × | ✓ | | | | | | 5 | 1 | 3 | 4 | 1 | 44.42 |
| **Others** | | | | | | | | | | | | | | | |
| Mapillary Vistas [96] | 2017 | 25K | × | | | ✓ | | | | 66 | 5 | 3 | 3 | 2 | 68.63 |
| TT 100K [15] | 2016 | 100K | × | | ✓ | | | | | 45 | 2 | 2 | 2 | 10 | 61.99 |
| Pascal3D+ [97] | 2014 | 30,899 | × | ✓ | ✓ | ✓ | | | | 12 | 2 | 2 | 1 | - | 58.07 |
| WildDash [98] | 2018 | 1,800 | × | | | ✓ | | | | 28 | 2 | 2 | 7 | 1 | 50.02 |
| TorontoCity [99] | 2016 | 56K | × | | | ✓ | | | | 4 | 2 | 2 | 1 | 1 | 45.82 |
| DAWN [100] | 2020 | 4,543 | × | ✓ | | | | | | 5 | 6 | 3 | 3 | - | 43.99 |
| RAD [101] | 2019 | 60 | × | | | ✓ | | | | 19 | 1 | 1 | 1 | 1 | 37.86 |
| STCrowd [102] | 2022 | 10,891 | ✓ | ✓ | ✓ | | | ✓ | | 1 | 3 | 1 | 1 | 1 | 35.19 |

The data recorded from Boston and Singapore cover varied driving behaviors and urban layouts, enhancing generalizability. Its six cameras provide a comprehensive perspective of the surrounding environment, making them widely utilized in multi-view object detection tasks. Yet, the dataset's coverage of rare driving cases, such as accidents, could be expanded to better assess algorithm performance under challenging conditions.

***Waymo.*** The Waymo Open Dataset [23], introduced in 2019, significantly influences research and advancement in autonomous driving. Waymo provides an extensive size of multimodal sensory data with high-quality annotations compared to others. Key contributions of the Waymo dataset include its comprehensive coverage of driving conditions and geographics, which are pivotal for the robustness and generability of different tasks. The real-world applicability of Waymo's data is enhanced by its diversity, although exploring its limitations in specific adverse conditions could provide deeper insights into areas needing improvement.

***BDD100K.*** BDD100K [24] dataset is renowned for its size and diversity. It comprises 100,000 videos, each about 40 seconds in duration. Furthermore, it includes different times of day and weather conditions, offering a solid foundation for testing and improving the robustness of algorithms. Meanwhile, it provides various annotated labels for object detection, tracking, semantic segmentation, and lane detection. However, the varying quality of video annotations challenges the dataset's real-world applicability, highlighting the need for consistent, high-quality annotations to ensure algorithm effectiveness across all conditions.

***Argoverse 2.*** As a sequel to Argoverse 1 [19], Argoverse 2 [28] introduces more diversified and complex driving scenarios, presenting the largest autonomous driving taxonomy to date. It captures various real-world driving scenarios across six cities and varying conditions. Argoverse 2 supports a wide range of essential tasks, including but not limited to 3D object detection, semantic segmentation, and tracking. The limitation of Argoverse 2 in the real world is similar to [22] and [23],

considering more adversarial conditions and edge cases can extend its application field.

*2) V2X: **TUMTraf.*** The TUM Traffic Dataset family (*TUMTraf*) is a cutting-edge real-world dataset comprising 50,253 labeled frames (9,545 point clouds and 40,708 images) across five releases, capturing diverse traffic scenarios in Munich, Germany [60], [84]–[86]. It integrates multiple data modalities like RGB camera, event-based camera, LiDAR, GPS, and IMU. TUMTraf also provides perspectives from the infrastructure and vehicle, enabling research in cooperative perception. TUMTraf stands out for including edge cases like accidents, near-miss events, and traffic violations, providing a great resource to improve perception systems.

***DAIR-V2X.*** DAIR-V2X [27] is a pioneering resource in the Vehicle-Infrastructure Cooperative Autonomous Driving, providing large-scale, multi-modality, multi-view real-world data. The dataset is designed to tackle challenges such as the temporal asynchrony between vehicle and infrastructure sensors and the data transmission costs involved in such cooperative systems. DAIR-V2X has been instrumental in advancing vehicle-infrastructure cooperation, setting benchmarks for addressing V2X perception tasks.

*3) Drone: **UAVDT.*** The UAVDT [94] dataset consists of 80,000 accurately annotated frames with up to 14 kinds of attributes, such as weather conditions, flying attitude, camera view, vehicle category, and occlusion levels. The dataset focuses on UAV-based object detection and tracking in urban environments. Moreover, the UAVDT benchmark includes high-density scenes with small objects and significant camera motion, all challenging for the current state-of-the-art methods.

***DroneVehicle.*** DroneVehicle [95] proposes a large-scale drone-based dataset, which provides 28,439 RGB-Infrared image pairs to address object detection, especially under low-illumination conditions. Furthermore, it covers a variety of scenarios, such as urban roads, residential areas, and parking lots. This dataset is a significant step forward in developing autonomous driving technologies due to its unique drone perspective across a broad range of conditions.

*4) Others: **Pascal3D+.*** Pasacal3D+ [97] is an extension of the *PASCAL VOC 2022* [103], overcoming the limitations of previous datasets by providing a richer and more varied set of annotations for images. Pasacal3D+ augments 12 rigid object categories, such as cars, buses, and bicycles, with 3D pose annotations and adds more images from ImageNet [104], resulting in high variability. However, Pascal3D+ only focuses on rigid objects and may not fully address the dynamic driving environments where pedestrians and other non-rigid objects are present.

***Mapillary Vistas.*** Mapillary Vistas dataset, proposed by [96] in 2017, particularly aims at semantic segmentation of street scenes. The 25,000 images in the dataset are labeled with 66 object categories and include instance-specific annotations for 37 classes. It contains images from diverse weather, time of day, and geometric locations, which helps mitigate the bias towards specific regions or conditions.

## B. Prediction, Planning, and Control Datasets

Prediction, planning, and control datasets serve as the foundation for facilitating the development of driving systems. These datasets are critical for forecasting traffic dynamics, pedestrian movements, and other essential factors that influence driving decisions. Hence, we demonstrate in detail several high-impact datasets related to these tasks according to the data size, modalities, and citation number. We summarize these datasets into task-specific and multi-task groups.

*1) Task-Specific Datasets: **highD.*** The drone-based highD [18] dataset provides a large-scale collection of naturalistic vehicle trajectories on German highways, containing post-processed trajectories of 110K cars and trucks. It is designed to address the shortcomings of traditional measurement techniques in scenario-based safety validation, which often miss capturing authentic road user behaviors or lack comprehensive, high-quality data. But highD is recorded under ideal weather conditions, limiting its application under adversarial weather conditions.

***PIE.*** The Pedestrian Intention Estimation (PIE) dataset proposed by [21] represents a significant advancement in understanding pedestrian behaviors in urban environments. It encompasses over 6 hours of driving footage recorded in downtown Toronto under various lighting conditions. PIE offers rich annotations for perception and visual reasoning, including bounding boxes with occlusion flags, crossing intention confidence, and text labels for pedestrian actions.

***Argoverse.*** Argoverse [19] is a crucial dataset in 3D object tracking and motion forecasting. Argoverse provides 360° images from 7 cameras, forward-facing stereo imagery, and LiDAR point clouds. The recorded data covers over 300K extracted vehicle trajectories from 290km of mapped lanes. With the assistance of rich sensor data and semantic maps, Argoverse is pivotal in advancing research and development in prediction systems. Real-world implementation has verified the effectiveness of Argoverse in urban environments, though its performance in different geographical areas may be constrained.

***nuPlan.*** nuPlan [113] is the world's first closed-loop machine learning-based planning benchmark in autonomous driving. This multimodal dataset comprises around 1,500 hours of human driving data from four cities across America and Asia. It features different traffic patterns, such as merges, lane changes, interactions with cyclists and pedestrians, and driving in construction zones. These characters of the nuPlan dataset take into account the dynamic and interactive nature of actual driving, allowing for a more realistic evaluation. Its real-world impact is seen in developing more adaptive and context-aware planning systems.

***exiD.*** The exiD [117] trajectory dataset, presented in 2022, is a pivotal contribution to the highly interactive highway scenarios. It takes advantage of drones to record traffic without occlusion, minimizing the influence on traffic and ensuring high data quality and efficiency. This drone-based dataset surpasses previous datasets regarding the diversity of interactions captured, especially those involving lane changes at highway entries and exits. Extending the data from weather conditions

TABLE III
PREDICTION, PLANNING, AND CONTROL DATASETS. WE DEMONSTRATE SEVERAL CRUCIAL DATASETS RELATED TO PREDICTION, PLANNING, AND CONTROL. BP: BEHAVIOR PREDICTION, DAP: DRIVER'S ATTENTION PREDICTION, IP: INTENTION PREDICTION, MP: MOTION PREDICTION, TP: TRAJECTORY PREDICTION, MPlan: MOTION PLANNING, DM: DECISION-MAKING, OT: OBJECT TRACKING, QA: QUESTION-ANSWERING, DBP: DRIVER BEHAVIOR RECOGNITION

| Dataset | Year | Sensing domain | Size | Tasks | Weather conditions | Time of day | Scenario conditions |
|---|---|---|---|---|---|---|---|
| Task-Specific | | | | | | | |
| Brain4Cars [105] | 2015 | others | 2M frames | Maneuver Anticipation | | | |
| JAAD [16] | 2017 | onboard | 75K frames | pedestrian IT | sunny, rainy, cloudy, snowy | day, afternoon, night | urban |
| Dr(eye)ve [59] | 2018 | onboard | 500K frames | DAP | sunny, rainy, cloudy | day, night | urban, countryside, highway |
| highD [18] | 2018 | drone | 45K km distance | TP | sunny | 8 am to 5 pm | highway |
| PIE [21] | 2019 | onboard | 293K frames | pedestrian IT | sunny, overcast | day | urban |
| USyd [106] | 2019 | onboard | 24K trajectories | driver IT | | | 5 intersections |
| Argoverse [19] | 2019 | onboard | 300K trajectories | TP | variety | variety | urban |
| Drive&Act [107] | 2019 | others | 9.6M images | DBR | | | |
| DbNet [108] | 2019 | onboard | 100 km distance | DBR | various road types | | |
| D²CAV [109] | 2020 | onboard | | behavioral strategy | | | |
| inD [110] | 2020 | drone | 11.5K trajectories | road user prediction | sunny | day | 4 urban intersections |
| PePscenes [111] | 2020 | onboard | 719 frames | pedestrian BP | | | |
| openDD [112] | 2020 | drone | 84,774 trajectories | pedestrian BP | | | 7 roundabouts |
| nuPlan [113] | 2021 | drone | 1.5K hours data | MPlan | | | 4 cities |
| DriPE [114] | 2021 | others | 10K pictures | DBR | | daytime | |
| Speak2label [115] | 2021 | others | 586 videos | DAP | sunny, cloudy | daytime | |
| CoCAtt [116] | 2022 | onboard | 11.9 hours | DAP | | | countryside |
| exiD [117] | 2022 | drone | 16 hours data | TP | sunny | daytime | 7 locations on highway |
| MONA [118] | 2022 | drone | 702K trajectories | TP | sunny, overcast, rain | 8 am to 5 pm | urban |
| Occ3D-nuScenes [119] | 2024 | onboard | 40K frames | occupancy prediction | | | |
| Occ3D-Waymo [119] | 2024 | onboard | 200K frames | occupancy prediction | | | |
| Multi-Task | | | | | | | |
| HDD [120] | 2018 | onboard | 104 hours data | driver behavior, causal reasoning | | | suburban, urban, highway |
| INTERACTION [20] | 2019 | drone, V2X | 110K trajectories | MPlan and MP, DM | | | (un)signalized intersection |
| BLVD [121] | 2019 | onboard | 120K frames | 4D OT, 5D event recognition | | day, night | urban, highway |
| rounD [122] | 2019 | drone | 13,746 road users | scenario classification, BP | sunny | daytime | (sub-)urban |
| PREVENTION [123] | 2019 | onboard | 356 mins video | IP, TP | | | highway, urban areas |
| DrivingStereo [124] | 2019 | onboard | 180K images | trajectory planning | sunny, rainy, cloudy foggy, dusky | daytime | suburban, urban, highway elevated road, country road |
| Lyft Level 5 [125] | 2021 | drone | 1.1K hours data | MPlan, MP | | | suburban |
| LOKI [126] | 2021 | onboard | 644 scenarios | TP, BP | variety | variety | (sub-)urban |
| SceNDD [127] | 2022 | onboard | 68 driving scenes | MPlan, MP | | | urban |
| DeepAccident [12] | 2023 | V2X | 57K frames | MP, accident prediction | sunny, rainy, cloudy, wet | noon, sunset, night | synthetic |
| V2X-Seq (forecasting) [56] | 2023 | V2X | 50K scenarios | online/offline VIC TP | | | 28 urban intersections |

and nighttime could be an improvement direction for this dataset.

**MONA.** The Munich Motion Dataset of Natural Driving (MONA) [118] is an extensive dataset, with 702K trajectories from 130 hours of videos, covering urban roads with multiple lanes, an inner-city highway stretch, and their transitions. This dataset boasts an average overall position accuracy of 0.51 meters, which exhibits the quality of the data collected using highly accurate localization and LiDAR sensors. However, only recording the data in a particular city may limit its generalizability to other geographic locations.

*2) Multi-Task Datasets: INTERACTION.* The INTERACTION [20] dataset is a versatile platform that offers diverse, complex, and critical driving scenes, along with a comprehensive semantic map, making it suitable for a multitude of tasks, such as motion prediction, imitation learning, and validation of decision and planning. Its inclusion of different countries and continents further improves the robustness of analyzing the driving behavior across different cultures. A potential shortcoming of the dataset is that the impact of environmental conditions has not been explicitly addressed.

**rounD.** The rounD dataset presented by [122] is pivotal for scenario classification, road user behavior prediction, and driver modeling. It provides a large number of road user trajectories at roundabouts. The dataset utilizes a drone equipped with a 4K resolution camera to collect over six hours of video, recording more than 13K road users. The broad recorded traffic situations and the high-quality recordings make rounD an essential dataset in autonomous driving, facilitating the study

of naturalistic driving behaviors in public traffic. Yet, similar to all datasets only collected under clear weather conditions, developing models capable of performing under challenging situations could be restricted.

**Lyft Level 5.** Lyft Level 5 [125] represents one of the most extensive autonomous driving datasets for motion prediction to date, with over 1,000 hours of data. It encompasses 17,000 25-second long scenes, a high-definition semantic map with over 15,000 human annotations, 8,500 lane segments, and a high-resolution aerial image of the area. It supports multiple tasks like motion forecasting, motion planning, and simulation. The numerous multimodal data with detailed annotations make it a vital benchmark for prediction and planning. The dataset shows limitations in accurately representing scenes with uncommon traffic conditions or rare pedestrian behaviors.

**LOKI.** LOKI [126], standing for Long Term and Key Intentions, is an essential dataset in multi-agent trajectory prediction and intention prediction. LOKI tackles a crucial gap in intelligent and safety-critical systems by proposing large-scale, diverse data for heterogeneous traffic agents, including pedestrians and vehicles. This dataset makes a multidimensional view of traffic scenarios available by utilizing camera images with corresponding LiDAR point clouds, making it a highly flexible resource for the community.

**DeepAccident.** The synthetic dataset, DeepAccident [12] is the first work that provides direct and explainable safety evaluation metrics for autonomous vehicles. The extensive dataset with 57K annotated frames and 285K annotated samples supports end-to-end motion and accident prediction, which is

vital in avoiding collisions and ensuring safety. Moreover, this multimodal dataset is versatile for various V2X-based perception tasks, such as 3D object detection, tracking, and BEV semantic segmentation. The various environmental conditions also enhance the generalizability of this dataset. Due to the domain adaptation issue, the performance of models trained on DeepAccident in real driving situations needs further study.

### C. End-to-end Datasets

End-to-end has become a growing trend as an alternative to modular-based architecture in autonomous driving [73]. Several versatile datasets (like nuScenes [22] and Waymo [23]) or simulators like *CARLA* [128] provide the opportunity to develop end-to-end autonomous driving. Meanwhile, some works present datasets that are especially for end-to-end driving.

**DDD17.** The DDD17 [129] dataset is notable for using event-based cameras. The dataset provides a concurrent stream of standard active pixel sensor (APS) images and dynamic vision sensors (DVS) temporal contrast events, offering a unique blend of visual data. Additionally, DDD17 captures diverse driving scenarios, including highway and city driving and various weather conditions, thus providing exhaustive and realistic data for training and testing end-to-end autonomous driving algorithms.

### VI. ANNOTATIONS PROCESS

The success and reliability of AD algorithms rely not only on the numerous data but also on high-quality annotations. In this section, we first explain the methodology for annotating data VI-A. Additionally, we analyze the most important aspects for ensuring annotation quality VI-B.

### A. Annotation Generation

Different AD tasks require specific types of annotation. For example, object detection requires bounding box labels of instance, segmentation is based on pixel- or point-level annotations, and continually labeled trajectory is critical for trajectory prediction. On the other hand, as shown in Fig. 11, the annotation pipeline can be categorized into three types: manual annotation, semi-automatic annotation, and fully automatic annotation. In this section, we detail the labeling approaches for different annotation types.

**Annotate 2D/3D Bounding Boxes.** The quality of the bounding box annotations directly impacts the effectiveness and robustness of the perception system (like object detection) of autonomous vehicles in real-world scenarios. The annotation process generally involves labeling images with rectangular boxes or point clouds with cuboids to precisely encompass the objects of interest.

*Labelme* [130] is a prior tool focusing on labeling images for object detection. However, generating bounding boxes by professional annotators faces the same issue as manual segmentation annotation. Wang et al. [131] presented a semi-automatic video labeling tool based on the open-source video annotation system *VATIC* [132]. Manikandan et al. [133] propose another automatic video annotation tool for AD scenes.



(a) Manual annotation

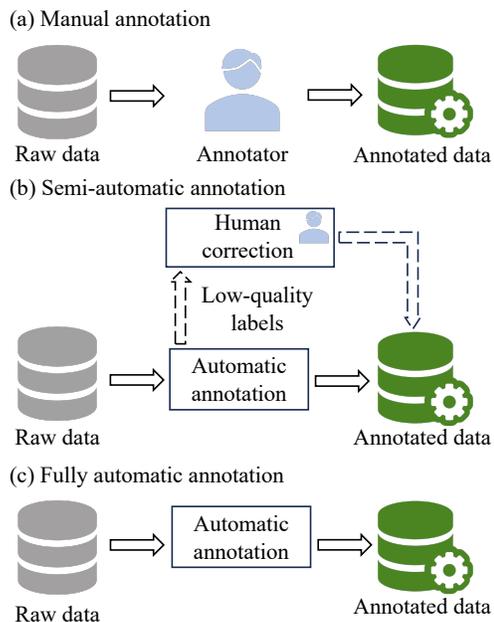(b) Semi-automatic annotation

(c) Fully automatic annotation

Fig. 11. Annotation pipelines. We demonstrate (a) Manual annotation: The professional annotators label the raw data using annotation tools. (b) Semi-automatic annotation: After generating annotations using an automatic annotation algorithm, the low-quality labels are refined by annotators. (c) Fully automatic annotation: The framework annotates data without human correction.

Process bounding box annotations in the nighttime is more challenging than daytime annotation. Schorkhuber et al. [134] introduced a semi-automatic approach leveraging the trajectory to solve this problem.

In contrast to 2D annotations, 3D bounding boxes contain richer spatial information, such as accurate location, the object's width, length, height, and orientation in space. Hence, labeling high-quality 3D annotations requires a more sophisticated framework. Meng et al. [135] applied a two-stage weakly supervised learning framework using human-in-the-loop to label LiDAR point clouds. *ViT-WSS3D* [136] generated pseudo-bounding boxes by modeling global interactions between LiDAR points and corresponding weak labels. *Apolloscape* [17] employed a labeling pipeline similar to [137], which consists of a 3D labeling and a 2D labeling branch, to handle static background/objects and moving objects, respectively. *3D BAT* [138] developed an annotation toolbox to assist in obtaining 2D and 3D labels in semi-automatic labeling.

**Annotate Segmentation Data.** The target of annotating segmentation data is to assign a label to each pixel in an image or each point in a LiDAR frame to indicate which object or region it belongs to. After labeling, all pixels belonging to an object are annotated with the same class. For the manual annotation process, the annotator first draws boundaries around an object and then fills in the area or paints over the pixels directly. However, generating pixel/point-level annotations in this way is costly and inefficient.

Many studies have proposed fully or semi-automatic annotation methods to improve annotation efficiency. Barnes et al. [139] presented a fully automatic annotation ap-

proach based on weakly supervised learning to segment proposed drivable paths in images. The semi-automatic annotation method [140] utilizes the objectness priors to generate segmentation masks. After that, [141] offered a semi-automatic method considering 20 classes. *Polygon-RNN++* [142] presents an interactive segmentation annotation tool following the idea of [143]. Instead of using image information to generate pixel-level labels, [137] explores transferring 3D information into 2D image domains to generate semantic segmentation annotations. For labeling 3D data, [144] proposes an image-assisted annotation pipeline. Luo et al. [145] leverage active learning to select a few points and to form a minimal training set to avoid labeling the whole point cloud scenes. Liu et al. [146] introduce an efficient labeling framework with semi/weakly supervised learning to label outdoor point clouds. **Annotate Trajectories.** A trajectory is essentially a series of points that map the path of an object over time, reflecting the spatial and temporal information. Labeling trajectory data for AD is a process that entails annotating the path or movement patterns of various entities within a driving environment, such as vehicles, pedestrians, and cyclists. Usually, the annotating process relies on object detection and tracking results.

As one of the prior works in trajectory annotation, [147] online generates actions for maneuvers and is annotated into the trajectory. The annotation approach [148] consists of a crowd-sourcing step followed by a precise process of expert aggregation. Jarl et al. [149] develop an active learning framework to annotate driving trajectory. Precisely anticipating movement patterns of pedestrians is critical for driving safety. Styles et al. [150] introduce a scalable machine annotation scheme for pedestrian trajectory annotations without human effort.

**Annotate on Synthetic Data.** Due to the expensive and time-consuming manual annotations on real-world data, synthetic data generated by computer graphics and simulators provide an alternative to address this issue. Since the data generation process is controllable and the attributes of each object in the scene (like position, size, and movement) are known, synthetic data can be automatically and accurately annotated.

The generated synthetic scenarios are designed to mimic real-world conditions, including multiple objects, various landscapes, weather conditions, and lighting variations. Selecting a suitable simulation tool is crucial to achieve this goal. *Torcs* [151] and *DeepDriving* [152] are two prior works for simulating autonomous vehicles while lacking multi-modality information and other types of objects like pedestrians. Recently, the open source simulators, such as *CARLA* [128], *SUMO* [153], and *AirSim* [154], have been widely used in data generation. They provide transparent platforms where researchers and developers can freely access and modify the source code, tailoring the simulator to specific requirements. In contrast, because of non-open sources, commercial simulation tools like *NVIDIA's Drive Constellation* [155] can make it difficult for users to create special driving environments. Rather than professional simulators, the game engines, including *Grand Theft Auto 5 (GTA5)* and *Unity* [156], are also popular for creating synthetic autonomous driving data.

Specifically, some researchers utilize the GTA5 game en-

gine to build datasets [4] and [5]. Krahenbuhl et al. [157] present a real-time system based on multiple games to generate annotations for various AD tasks. Instead of applying game videos, *SHIFT* [10], *CAOS* [7], *FedBEVT* [158], and *V2XSet* [11] are created based on the CARLA simulator. Compared to [11], *V2X-Sim* [89] studies employing multiple simulators [128], [153] to generate dataset for V2X perception tasks. *CODD* [159] further exploits using [128] to generate 3D LiDAR point clouds for cooperative driving. Other works [2], [3], [9], [160] leverage the Unity development platform to generate synthetic datasets.

### B. Quality of Annotations

Supervised learning-based autonomous driving (AD) algorithms rely heavily on extensive, well-labeled datasets. High-quality datasets ensure that these systems can accurately perceive and interpret complex driving environments, enhancing safety and reliability on the road. This in turn fosters user trust, a crucial factor for the widespread adoption of autonomous vehicles. Conversely, poor dataset quality can result in system errors and safety risks, undermining user confidence, hindering acceptance, and failing to meet the criteria for trustworthy AI as discussed by [161]. Therefore, ensuring the quality of annotations is fundamental for improving accuracy while driving in complex real-world environments. It is even more important to fine-tune the data quality and use active learning methods for dataset curation to get robust performance results on a test set than fine-tuning the model architecture [162].
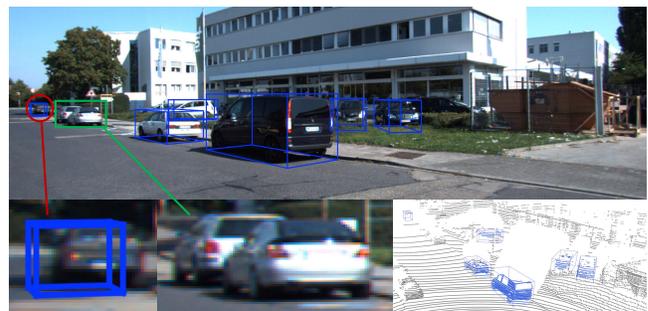


Fig. 12. Mislabeling example of KITTI [13] dataset. We show the ground truth in blue. The bounding box of a car (in red circle) is not precise. Two cars (in green cube) are not annotated, although sensors obviously capture them.

According to the study [163], the annotation quality is affected by several factors, such as consistency, correctness, precision, and validation. Consistency is the foremost criterion in evaluating annotation quality. It involves maintaining uniformity across the entire dataset and is crucial for avoiding confusion in models trained on this data. For example, if a particular type of vehicle is labeled as a car, it should be consistently annotated the same in all other instances. Annotation precision is another vital indicator, which refers to whether the labels match the actual state of the objects or scenarios. In contrast, correctness highlights that annotated data are appropriate and relevant for the purpose of

the dataset and annotation guidelines. After annotation, it is essential to validate the annotated data to ensure its accuracy and completion. The process can be done through manual review by experts or algorithms. Validation helps effectively prevent issues in datasets before they infect the performance of autonomous vehicles, decreasing potential safety risks. Inel et al. [164] presented a data-agnostic validation method for expert annotated datasets.

A failure case of annotation from KITTI [13] is shown in Fig. 12. We illustrate the ground truth bounding boxes (blue) in the corresponding image and LiDAR point cloud. On the left side of the image, the annotation of a car (circled in red) is inaccurate because it does not contain the whole object car. Additionally, two cars (highlighted by the green rectangle) are not annotated, even though the camera and LiDAR capture them clearly. Furthermore, datasets like the *IPS300+* [165] contain many labeled objects within a frame (319.84 labels per frame), but on the other hand, the annotations are of bad quality. Many large datasets like the *Pandaset* [166], *Oxford* [167], *CADC* [168], nuScenes [22], and *Lyft Level 5* [22] were labeled by specialized labeling companies like *Scale AI* and provide high-quality annotations. Labeling the nuScenes dataset took about 7,937 hours and cost 100k USD. Another way to label datasets is to use custom labeling tools like *3D BAT* [60], [138], which was used to create the *TUMTraf* dataset. The Waymo and KITTI datasets were both labeled with custom labeling tools. *V2V4Real* [87] has used the *SUSTechPoints* [169] labeling tool to generate the dataset.

## VII. DATA ANALYSIS

In this section, we systematically analyze datasets from different perspectives in detail, such as the distribution of data around the world (VII-A), chronological trend VII-B, and the data distribution VII-C.
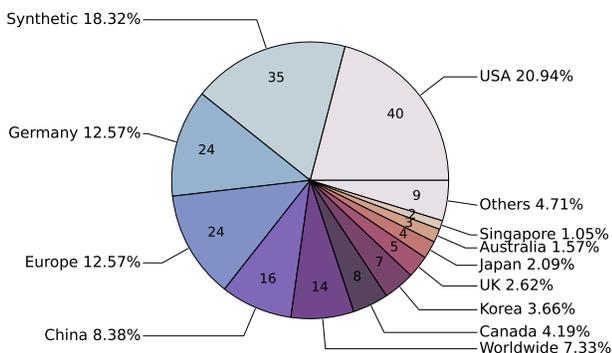
### A. Worldwide Distribution



Fig. 13. The distribution of datasets around the world. This figure illustrates the distribution of data collection locations of the datasets.

We demonstrate an overview of the global distribution of 191 autonomous driving datasets in Fig. 13. The chart indicates that the USA is at the forefront with 40 datasets (21% share), underscoring its leadership in the autonomous

driving domain. Germany accounts for 24 datasets, and China closely follows with 16 datasets. On the other hand, developed countries, including Canada, Korea, the UK, Japan, and Singapore, cover the rest smaller segments. Although 11 datasets are collected worldwide and 24 are from the European region (except Germany), all these countries or regions are considered high-income areas. The dominance of the USA, western Europe, and East Asia reflects an extremely unbalanced development of autonomous driving around the world.

Specifically, one of the most classic datasets, KITTI [13], is collected in the metropolitan area of Karlsruhe, Germany. In comparison, both Waymo [23] and Argoverse 2 [28] datasets are collected from a wide variety of places, including six different cities in the USA individually. Apolloscapes [17] and DAIR-V2X [27] are recorded in China. Instead of collecting data solely within a country, nuScenes [22] is established based on the data from America (Boston) and Singapore, both known for complex and highly challenging traffic situations. Other well-known autonomous driving datasets [18], [21], [24], [88], [95], [118] are collected in those previously mentioned countries. It is noted that thanks to the geographic diversity, [22], [23] have been widely used in transfer learning to verify the generalizability of autonomous driving algorithms.

Moreover, autonomous driving faces unique challenges across different geographical regions. However, relying solely on data from a single source can introduce bias that could result in autonomous vehicles' failure to perform under varied or unseen regions and cases. For instance, the diversity and quantity of electric scooters in China far exceed those in Germany, meaning an algorithm trained exclusively on German data may struggle to accurately recognize objects in China. Hence, Recording data from different continents and countries can assist in addressing unique challenges caused by geographical locations. This diverse regional distribution enhances the robustness of the collected data and highlights the international efforts and collaborations in the research community and industry.

Furthermore, 35 synthetic datasets generated by simulators like CARLA [128] take up 18.32% percent. Due to the limitation of recording from real-world driving environments, these synthetic datasets overcome such drawbacks and are critical for exploiting more robust and reliable driving systems. However, the domain adaptation from synthetic data to real-world data is still a challenging research topic, which limits a broader application of synthetic data and relevant simulators.

### B. Chronological Trends in Perception Datasets

In Fig. 10, we introduce a chronological overview of perception datasets with the top 50 impact scores from 2009 to 2024 (until the writing of this paper). The datasets are color-coded according to their sourcing domain, and synthetic datasets are marked with an external red outline, clearly illustrating the progress toward the diverse data collection strategy. A noticeable trend shows the increase in the number and variety of datasets over the years, indicating the requirement for high-quality datasets with the growing advancements in the field of autonomous driving.

In general, most of the datasets provide a perception perspective from the sensors equipped on the ego vehicle (onboard) because of the importance of the capability of the autonomous vehicle to efficiently and precisely precept the surroundings. On the other hand, due to the high-cost real-world data, some researchers propose high-influence synthetic datasets like VirtualKITTI [3] (2016) to alleviate the reliance on real data. Facilitated by the effectiveness of simulators, there are many novel synthetic datasets [7] [10] published in recent years. In the timeline, V2X datasets like DAIR-V2X [27] and TUMTraf ( [60], [84]–[86]) also exhibit a trend toward cooperative driving systems. Furthermore, because of the non-occlusion perspective provided by UAV, drone-based datasets, such as UAVDT [94] published in 2018, take a crucial position in advancing perception systems.
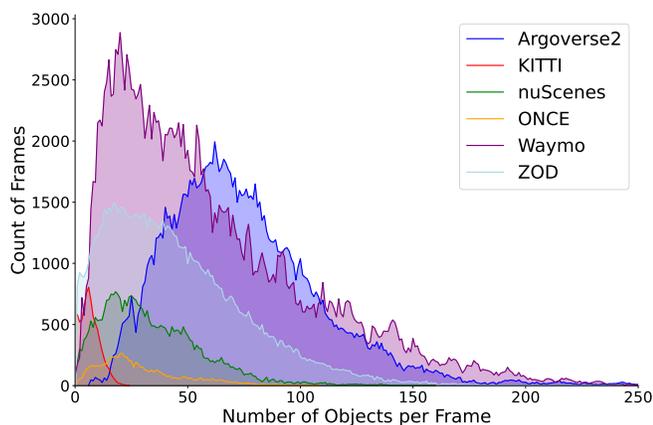
### C. Data Distribution



Fig. 14. Comparison of the distribution of the number of objects per frame across several datasets: Argoverse 2 [28], KITTI [13], nuScenes [22], ONCE [30], Waymo [23], and ZOD [31]. The horizontal axis quantifies the number of objects detected in a single frame, while the vertical axis represents the count of frames containing that number of objects.
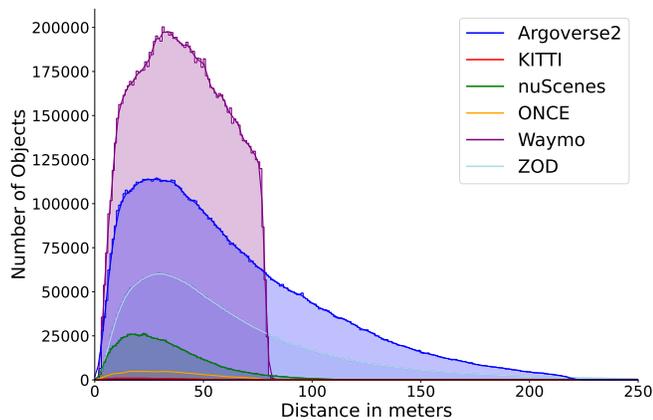


Fig. 15. Comparison of the distribution of the number of objects detected at various distances across several datasets: Argoverse 2 [28], KITTI [13], nuScenes [22], ONCE [30], Waymo [23], and ZOD [31]. The horizontal axis measures the distance from the ego vehicle in meters, and the vertical axis quantifies the number of objects detected at that distance.

We introduce an insight into the number of objects per frame for these datasets in Fig. 14. Notably, Waymo [23] exhibits an extreme number of frames with less than 50 objects while maintaining a broad presence across the chart, illustrating a wide range of scenarios from low to high object density per frame. Contrastingly, KITTI [13] shows a more constrained distribution and limited data size. Argoverse 2 [28] features a substantial number of frames with a higher object count—its peak appears around 70, which indicates its complex environmental representations in general. For ONCE [30], its density of objects evenly distributes in the supported perception range. Datasets like nuScenes [22] and ZOD [31] demonstrate similar curves with a quick rise and slow decline, implying a moderate level of environmental complexity with a decent variability of object counts per frame.

Beyond the number of objects in a scene, the object distribution based on the distance to the ego vehicle is another essential point for revealing a dataset's variety and significant differences, illustrated in Fig. 15. The Waymo dataset demonstrates numerous labeled objects in near-field to mid-field scenarios. In contrast, Argoverse 2 and ZOD show a wider detection range, with some frames even including bounding boxes out of 200 meters. The curve of the nuScenes means it is particularly rich in objects in a shorter range, which is typical for urban driving scenarios. Nevertheless, as the distance increases, the nuScenes dataset quickly tappers off for the number of objects with annotations. The ONCE dataset covers a more even distribution of objects across distances, while the KITTI dataset focuses more on close-range objects.

### D. Influence of Adversarial Environmental Conditions.

TABLE IV
3D OBJECT DETECTION PERFORMANCE UNDER ADVERSARIAL ENVIRONMENTAL CONDITIONS ON THE NUSCENES [22] DATASET. WE REPORT MAP (%) AS THE EVALUATION METRIC. 'L' AND 'C' REFER TO LIDAR AND CAMERAS, RESPECTIVELY. THE TERM 'NORMAL' INDICATES THAT THE MODEL WAS TESTED ACROSS THE ENTIRE VALIDATION SET, ENCOMPASSING ALL ENVIRONMENTAL CONDITIONS.

| Method | Modality | Adversarial Conditions | mAP ↑ |
|---|---|---|---|
| VoxelNext [170] | L | normal | **60.5** |
| | | nighttime | 32.4 |
| | | rainy | 58.7 |
| UVTR [171] | L | normal | **60.8** |
| | | nighttime | 33.9 |
| | | rainy | 56.5 |
| Transfusion [172] | L+C | normal | **68.9** |
| | | nighttime | 37.5 |
| | | rainy | 65.8 |

We further study the influence of adversarial environmental conditions (low lightning and rain) on the performance of 3D object detectors in the autonomous driving system. The experiment results are reported in Tab. IV. We utilize the nuScenes [22] dataset and choose three state-of-the-art methods, *VoxelNext* [170], *UTVR* [171], and *Transfusion* [172]. For a fair comparison, we directly utilize the pre-trained checkpoints provided in the repository of each approach. The rainy and nighttime data are also manually collected from the

TABLE V
VLM AUTONOMOUS DRIVING DATASET. OT: OBJECT TRACKING, MOT: MULTI-OBJECT TRACKING, QA: QUESTION-ANSWERING, DM:
DECISION-MAKING, IP: INTENTION PREDICTION, VR: VISUAL REASONING, SR: SPATIAL REASONING

| Dataset | Year | Size | Temporal | Sensing domain | Tasks | Real/Synthetic |
|---|---|---|---|---|---|---|
| BDD-X [175] | 2018 | 8.4M | ✓ | onboard | reasoning, planning | real |
| Cityscapes-Ref [176] | 2018 | 5,000 stereo videos | ✓ | onboard | object referring | real |
| TOUCHDOWN [177] | 2019 | 9,326 examples | ✓ | onboard | reasoning, navigation | real |
| Talk2Car [178] | 2019 | 11,959 commands, 850 videos | ✓ | onboard | object referring | real |
| BDD-OIA [179] | 2020 | 11,303 scenarios | ✗ | onboard | explainable decision-making | real |
| CityFlow-NL [180] | 2021 | 5,289 samples | ✓ | onboard | OT | real |
| CARLA-NAV [181] | 2022 | 83K | ✓ | onboard | navigation | real |
| NuPrompt [182] | 2023 | 34K frames, 35K prompts | ✓ | onboard | MOT | real |
| NuScenes-QA [183] | 2023 | 1K scenarios, 460K QA pairs | ✓ | onboard | visual QA | real |
| Refer-KITTI [57] | 2023 | 6,650 frames | ✓ | onboard | referring MOT | real |
| Driving LLMs [184] | 2023 | 10K driving situations, 160K QA pairs | ✗ | drone | visual QA | synthetic |
| DRAMA [185] | 2023 | 17,785 scenarios | ✓ | onboard | reasoning, visual QA | real |
| Rank2Tell [186] | 2023 | 116 scenarios | ✓ | onboard | importance level ranking | real |
| LamPilot [187] | 2023 | 4,900 samples | ✓ | others | planning | synthetic |
| LangAuto CARLAR [188] | 2023 | 64K data clips | ✓ | onboard | closed-loop driving | synthetic |
| NuScenes-MQA [189] | 2023 | 34K scenarios, 1.4M QA pairs | ✗ | onboard | visual QA | real |
| DriveMLM [190] | 2023 | 280 hours | ✓ | onboard | planning, control | synthetic |
| DriveLM-nuScenes [191] | 2023 | 4,871 frames | ✓ | onboard | end-to-end driving | real |
| DriveCARLA [191] | 2023 | 183,373 frames | ✓ | onboard | end-to-end driving | real |
| LiDAR-text [192] | 2023 | 420K 3D captioning data, 280K 3D grounding data | ✗ | onboard | 3D scene understanding | real |
| Talk2BEV [193] | 2023 | 1K BEV scenarios, 20K QA pairs | ✓ | onboard | DM, IP, VR, SR | real |
| NuScenes-MQA [194] | 2024 | 1.4M QA pairs, 34,149 scenarios | ✓ | onboard | visual QA | real |
| VLAAD [195] | 2024 | 10,379 scenarios | ✓ | onboard | visual QA, reasoning | real |

nuScenes validation set. All three methods demonstrate similar trends under different environmental conditions. Specifically, the detection accuracy declines significantly under low illumination conditions compared with evaluations of the whole validation set. Besides, all models exhibit slight decreases under the non-heavy rainy situations recorded by nuScenes. The reliability of detectors in the real world could be much worse once the rain is heavy. Therefore, taking image restoration into account is a promising way for camera or sensor fusion-based approaches to overcome these challenges [173], [174]. In conclusion, increasing the amount of data recorded on various types and intensities of weather conditions is critical for training a robust and reliable autonomous driving system.

## VIII. DISCUSSION AND FUTURE WORKS

With rapid technological development, powerful computation resources, and excellent artificial intelligent algorithms, many novel trends in next-generation autonomous driving datasets have occurred while proposing new challenges and requirements.

**End-to-End Driving Datasets.** Compared to the modular-designed autonomous driving pipeline, the end-to-end architecture simplifies the overall design process and reduces integration complexities. The success of *UniAD* [196] verifies the potential ability of end-to-end models. However, the number of datasets for end-to-end AD is limited [129]. Therefore, introducing datasets focusing on end-to-end driving is crucial for advancing autonomous vehicles. On the other hand, implementing an automatic labeling pipeline in a data engine can significantly facilitate the development of end-to-end driving frameworks and data [74].

**Potential Applications of AD Datasets** Future autonomous driving datasets should provide extensive real-world environmental and traffic data, supporting a wide range of applications beyond the ego-vehicle and basic vehicle-infrastructure cooperation. For instance, the interaction data between autonomous vehicles and intelligent infrastructure components can guide the advancement of Internet-of-Thing (IoT)-enabled devices like smart traffic signals. Moreover, the detailed insights into traffic patterns, congestion, and vehicle behavior across different times and conditions can facilitate urban planning, optimize the traffic flow, and enhance overall traffic management strategies.

**Introduce Language into AD Datasets.** Vision language models (VLMs) have recently achieved impressive advancement in many fields. Its inherent advantage in providing language information to vision tasks makes autonomous driving systems more explainable and reliable. The survey [197] highlights the prominent role of Multimodal Large Language Models in various AD tasks, such as perception [176], [182], motion planning [187], and motion control [193]. Autonomous driving datasets including language labels are shown in Tab. V. Overall, incorporating language into AD datasets is a trend of the future development of AD datasets.

**Data Generation via VLMs.** The powerful capability of VLMs can be used to generate data [198]–[201]. For example, *DriveGAN* [202] generated high-quality AD data by disentangling different components without supervision. Additionally, due to the capability of world models for comprehending driving environments, some works [203]–[205] explored world models to generate high-quality driving videos. *Drive-Dreamer* [204] as a pioneering work derived from real-world scenarios, addressing the limitation of gaming environments or simulated settings. The most recent advancements in text-video generation, exemplified by *Sora* [206], have opened a new avenue for generating autonomous driving data. With a brief description, Sora can generate high-quality synthetic data that is very close to real scenes. This capability significantly enhances dataset augmentation, especially by extending the data available for rare events like traffic accidents. The increased data availability provided by VLMs is poised to enhance the training and evaluation of autonomous driving

systems, potentially improving their safety and reliability.

**Domain Adaptation.** Domain adaptation is a critical challenge in developing autonomous vehicles [207], referring to the ability of a model trained on one dataset (the source domain) to perform stable on another dataset (the target domain). This challenge manifests in multiple aspects, such as diversity in driving conditions [208], sensor settings [209], or synthetic-to-real transform [210]. Therefore, a trend of developing the next-generation datasets is incorporating data from heterogeneous sources. First of all, datasets should not only cover a wide range of environmental conditions (various weather conditions and day-night illumination) but also include diverse geographic locations. Second, combining data from various sensor types is also pivotal to overcoming the domain adaptation issue. Another solution is taking into account blending high-quality synthetic data and real-world data in a balanced way to improve generalization.

**Uncertainty Issues in Autonomous Driving.** Usually, uncertainty is modeled in a probabilistic way in the field of machine learning. The uncertainty can be referred to two types: a) aleatoric uncertainty, where the uncertainty like noise stems from the data, and b) epistemic uncertainty refers to uncertainty caused by unawareness about the best model [211]. For autonomous driving, one of the major reasons resulting in uncertainty is the incompleteness of training data [212]. The insufficient training data can not wholly represent the driving environment, causing autonomous vehicles to problematically deal with rare cases during operation. Therefore, improving the diversity of the dataset for training reliable autonomous driving systems is crucial. Moreover, datasets include rare events and edge cases, allowing models to better understand and quantify the uncertainty associated with unexpected situations and safely handle them.

**Standardization of Data Creation.** Addressing the standardization of data creation is critical aspect for developing new datasets, as it directly impacts the accuracy, efficiency, and reliability of autonomous vehicle models. In general, data standardization refers to the attributes of the data, terminology and structure of the dataset, and data storage [213]. For the data attributes, unifying data formats across different sensor types and source domains is valuable to facilitate integrated processing and analysis. Developing comprehensive guidelines for dataset labeling is crucial to guarantee consistent, high-quality annotations, enhancing model training across diverse datasets and boosting performance reliability. Moreover, establishing standardized data storage and access protocols is vital, enabling seamless dataset sharing and integration across multiple sources, and fostering collaboration within the autonomous driving research community.

**Data Privacy.** The advancement of autonomous vehicles require a lot of data to guarantee the driving safety. However, the more the available data, the greater the concern that private data will be abused. For the earlier autonomous driving datasets like KITTI [13], there is no anonymization progress for the recorded data, especially for images, which leads to a potential risk for leak private information. With the introduction and refinement of relevant regulations in various countries, data anonymization has been widely adopted in recent datasets [22], [23], [84]. Nevertheless, the information provided by a large dataset is more than individual biographic features or vehicle plates. Professional institutes can indirectly extract external information by analyzing the existing information in the data. For instance, by analyzing the types of vehicles and the attire of pedestrians in the dataset, one can infer information about the infrastructure, basic construction, and other characteristics of the surrounding area in which the data were collected.

**Open Data Ecosystems.** The primary objective of open data ecosystems (ODE) for autonomous driving is to foster innovation, enhance transparency, and facilitate collaboration across governments, companies, and research communities. This is achieved through the free exchange of datasets, breaking down the traditional barriers that have restricted access to corporations and research institutions. By doing so, ODE empowers a wider range of innovators to participate in autonomous driving development, thereby boosting a more diverse and inclusive innovation ecosystem. Additionally, ODE establishes dynamic feedback loops where users can report issues, propose improvements, and contribute to enhancing datasets. Nevertheless, unrestricted access to data raises strong concerns about data security and privacy. Addressing these concerns necessitates carefully crafting and continuously refining relevant legal frameworks to safeguard sensitive information while prompting the growth of ODE. The balance is pivotal for maintaining the integrity of open data ecosystems and ensuring their sustainable development in autonomous driving.

## IX. Conclusion

In this paper, we exhaustively and systematically reviewed 265 existing autonomous driving datasets. We started with the sensor types and modalities, sensing domains, and tasks relevant to autonomous driving. We introduced a novel evaluation metric called impact score to validate the influence and importance of perception datasets. Our analysis covered the attributes and significance of high-impact datasets across perception, prediction, planning, control, and end-to-end driving tasks. Additionally, we investigated the annotation methodologies and the factors affecting annotation quality. We also analyzed datasets from chronological and geographical perspectives to understand current trends and conducted experiments demonstrating the crucial need for data diversity under adversarial conditions. With the study on data distribution, we offered a specific viewpoint for comprehending the variances across different datasets. Our findings underscore the essential role of diverse, high-quality datasets in shaping the future of autonomous driving. Looking forward, we highlighted key challenges and future directions for autonomous driving datasets, including the adoption of VLM, domain adaptation, uncertainty challenges, standardization, data privacy, and open data ecosystems. These areas represent promising avenues for future research and are pivotal for advancing autonomous driving technology, setting a clear roadmap for further innovation in this rapidly evolving field.

## Appendix

### Supplementary Tables of AD Datasets

TABLE VI

PART I OF AUTONOMOUS DRIVING DATASETS. SS: SEMANTIC SEGMENTATION, OD: OBJECT DETECTION, OT: OBJECT TRACKING, MOT: MULTI-OBJECT TRACKING

| Dataset | Year | Size | Temporal | Sensing domain | Tasks | Real/Synthetic |
|---|---|---|---|---|---|---|
| ETH Ped [214] | 2007 | 2,293 frames | ✓ | onboard | pedestrian detection | real |
| TUD-Brussels [215] | 2009 | 1,600 frames | ✗ | onboard | (3D) OD | real |
| Collective activity [216] | 2009 | 44 short videos | ✓ | others | human activity recognition | real |
| San Francisco Landmark [217] | 2011 | 150K panoramic images | ✗ | others | landmark identification | real |
| Daimler Stereo Ped [218] | 2011 | 21,790 frames | ✓ | onboard | pedestrian detection | real |
| BelgiumTS [219] | 2011 | 13K traffic sign annotations | ✗ | onboard | traffic sign detection | real |
| Stanford Tack [220] | 2011 | 14K tracks | ✓ | onboard | classification | real |
| TME Motorway [221] | 2012 | 30K frames | ✓ | onboard | OD, OT | real |
| MSLU [222] | 2013 | 36.8 km distances | ✓ | onboard | SLAM | real |
| SydneyUrbanObject [223] | 2013 | 588 object scans | ✗ | onboard | classification | real |
| Ground Truth SitXel [224] | 2013 | 78,500 frames | ✓ | onboard | stereo confidence | real |
| NYC3DCars [225] | 2013 | 2K images | ✗ | onboard | OD | real |
| AMUSE [226] | 2013 | 117,440 frames | ✓ | onboard | SLAM | real |
| Daimler Ped [227] | 2013 | 12,485 frames | ✓ | onboard | pedestrian path prediction | real |
| LISA [228] | 2014 | 6,610 frames | ✓ | onboard | traffic sign detection | real |
| Paris-rue-Madame [229] | 2014 | 643 objects | ✗ | onboard | OD, SS | real |
| TRANCOS [230] | 2015 | 1.2K images | ✗ | V2X | onboard number estimation | real |
| FlyingThings3D [231] | 2015 | 26,066 frames | ✓ | others | scene flow estimation | synthetic |
| Ua-detrac [232] | 2015 | 140K frames | ✓ | V2X | OD, MOT | real |
| NCLT [233] | 2015 | 34.9 hours | ✓ | onboard | odometry | real |
| CCSAD [234] | 2015 | 96K frames | ✓ | onboard | scene understanding | real |
| KAIST MPD [235] | 2015 | 95K color-thermal pair frames | ✗ | onboard | pedestrian detection | real |
| SAP [236] | 2016 | 19K frames | ✗ | drone | OD, OT | real |
| LostAndFound [237] | 2016 | 2,104 frames | ✓ | onboard | obstacle detection | real |
| UAH-Driveset [238] | 2016 | 500 mins | ✓ | onboard | lane detection, detection | real |
| CURE-TSR [239] | 2017 | 2.2M annotated images | ✗ | onboard | traffic sign detection | real |
| TuSimple [240] | 2017 | 6,408 frames | ✗ | onboard | lane detection, velocity estimation | real |
| TRoM [241] | 2017 | 712 frames | ✗ | onboard | road marking detection | real |
| NEXET [242] | 2017 | 91,190 frames | ✗ | onboard | OD | real |
| DIML [243] | 2017 | 470 videos | ✓ | onboard | lane detection | real |
| Bosch STL [244] | 2017 | 13,334 images | ✓ | onboard | traffic light detection and classification | real |
| Complex Urban [245] | 2017 | around 190km paths | ✓ | onboard | SLAM | real |
| DDD20 [246] | 2017 | 51 hours event frames | ✓ | onboard | end-to-end driving | real |
| PedX [247] | 2018 | 5K images | ✓ | onboard | pedestrian detection and tracking | real |
| LiVi-Set [248] | 2018 | 10K frames | ✓ | onboard | driving behavior prediction | real |
| Syncapes [249] | 2018 | 25K images | ✗ | onboard | OD, SS | synthetic |
| NVSEC [250] | 2018 | around 28 km distances | ✓ | others | SLAM | real |
| Aachen Day-Night [251] | 2018 | 4,328 images, 1.65M points | ✓ | onboard | visual localization | real |
| CADP [252] | 2018 | 1,416 scenes | ✗ | V2X | traffic accident analysis | real |
| TAF-BW [253] | 2018 | 2 scenarios | ✓ | V2X | MOT, V2X communication | real |
| comma2k19 [254] | 2018 | 2M images | ✗ | onboard | pose estimation, end-to-end driving | real |
| nighttime drive [255] | 2018 | 35K | ✗ | onboard | SS | real |
| VEIS [160] | 2018 | 61,305 frames | ✗ | onboard | OD, SS | synthetic |
| Paris-Lille-3D [256] | 2018 | 2,479 frames | ✗ | onboard | 3D SS, classification | real |
| NightOwis [257] | 2018 | 56K frames | ✓ | onboard | pedestrian detection, tracking | real |
| EuroCity Persons [258] | 2018 | 47,300 frames | ✗ | onboard | OD | real |
| RANUS [259] | 2018 | 4K frames | ✗ | onboard | SS, scene understanding | real |
| SynthCity [260] | 2019 | 367.9M points in 30 scans | ✗ | onboard | (3D) OD, (3D) SS | synthetic |
| D²-City [261] | 2019 | 700K annotated frames | ✓ | onboard | OD, MOT | real |
| Caltech Lanes [262] | 2019 | 1,224 frames | ✗ | onboard | lane detection | real |
| Mcity [263] | 2019 | 1,7500 frames | ✓ | onboard | SS | real |
| DET [264] | 2019 | 5,424 event-based camera images | ✗ | onboard | lane detection | real |
| PreSIL [8] | 2019 | 50K frames | ✗ | onboard | OD, 3D SS | synthetic |
| H3D [265] | 2019 | 27,721 frames | ✓ | onboard | (3D) OD, MOT | real |
| LLAMAS [266] | 2019 | 100K frames | ✓ | onboard | lane segmentation | real |
| MIT DriveSeg [267] | 2019 | 5K frames | ✗ | onboard | SS | real |
| Astyx [268] | 2019 | 500 radar frames | ✗ | onboard | 3D OD | real |
| UNDD [269] | 2019 | 7.2K frames | ✓ | onboard | SS | real |
| Boxy [270] | 2019 | 200K frames | ✓ | onboard | OD | real |
| RUGD [271] | 2019 | 37K frames | ✓ | others | SS | real |
| ApolloCar3D [272] | 2019 | 5,277 driving images | ✗ | onboard | 3D instance understanding | real |
| HEV [273] | 2019 | 230 video clips | ✓ | onboard | object localization | real |
| HAD [274] | 2019 | 5,675 video clips | ✓ | onboard | end-to-end driving | real |
| CARLA-100 [275] | 2019 | 100 hours driving | ✓ | onboard | path planning, behavior cloning | synthetic |
| Brno Urban [276] | 2019 | 375.7 km | ✓ | onboard | recognition | real |
| VERI-Wild [277] | 2019 | 416,314 images | ✓ | V2X | onboard re-identification | real |
| CityFlow [278] | 2019 | 200K bounding boxes | ✓ | V2X | OD, MOT, re-identification | real |
| VLMV [279] | 2020 | 306K frames | ✓ | V2X | lane merge | real |
| Small Obstacles [280] | 2020 | 3K frames | ✗ | onboard | small obstacle segmentation | real |
| Cirrus [281] | 2020 | 6,285 frames | ✓ | onboard | (3D) OD | real |
| KITTI InstanceMotSeg [282] | 2020 | 12,919 frames | ✓ | onboard | moving instance segmentation | real |
| A*3D [283] | 2020 | 39,179 point cloud frames | ✗ | onboard | (3D) OD | real |
| Toronto-3D [284] | 2020 | 4 scenarios | ✗ | onboard | 3D SS | real |
| MIT-AVT [285] | 2020 | 1.15M 10s video clips | ✓ | onboard | SS, anomaly detection | real |
| CADC [168] | 2020 | 56K | ✗ | onboard | (3D) OD, OT | real |
| SemanticPOSS [286] | 2020 | 2,988 point cloud frames | ✓ | onboard | 3D SS | synthetic |
| IDDA [287] | 2020 | 1M frames | ✗ | onboard | segmentation | synthetic |

TABLE VII
PART II OF AUTONOMOUS DRIVING DATASETS

| Dataset | Year | Size | Temporal | Sensing domain | Tasks | Real/Synthetic |
|---|---|---|---|---|---|---|
| CARRADA [288] | 2020 | 7,193 radar frames | ✓ | onboard | SS | real |
| Titan [289] | 2020 | 75,262 frames | ✓ | onboard | OD, action recognition | real |
| NightCity [290] | 2020 | 4,297 frames | ✗ | onboard | nighttime SS | real |
| PePScenes [111] | 2020 | 40K frames | ✓ | onboard | (3D) OD, pedestrian action prediction | real |
| DDAD [291] | 2020 | 21,200 frames | ✗ | onboard | depth estimation | real |
| MulRan [292] | 2020 | 41.2km paths | ✓ | onboard | place recognition | real |
| Oxford Radar RobotCar [167] | 2020 | 240K scans | ✓ | onboard | odometry | real |
| OTOH [125] | 2020 | 170K scenes | ✓ | drone | trajectory prediction, planning | real |
| DA4AD [293] | 2020 | 9 sequences | ✓ | onboard | visual localization | real |
| CPIS [92] | 2020 | 10K frames | ✗ | V2X | cooperative 3D OD | synthetic |
| EU LTD [294] | 2020 | around 37 hours | ✓ | onboard | odometry | real |
| Newer College [295] | 2020 | 290M points, 2300 seconds | ✓ | others | SLAM | real |
| CCD [296] | 2020 | 4.5K videos | ✓ | onboard | accident prediction | real |
| LIBRE [297] | 2020 | 40 frames | ✗ | onboard | LiDAR performance benchmark | real |
| Gated2Depth [298] | 2020 | 17,686 frames | ✓ | onboard | depth estimation | real |
| TCGR [299] | 2020 | 839,350 frames | ✓ | others | traffic control gesture recognition | real |
| DSEC [47] | 2021 | 53 sequences (3193 seconds in total) | ✓ | onboard | dynamic perception | real |
| 4Seasons [300] | 2021 | 350km recordings | ✓ | onboard | SLAM | real |
| PVDN [301] | 2021 | 59,746 frames | ✓ | onboard | nighttime OD, OT | real |
| ACDC [80] | 2021 | 4,006 images | ✗ | onboard | SS on adverse conditions | real |
| DRIV100 [302] | 2021 | 100K frames | ✗ | onboard | domain adaptation SS | real |
| NEOLIX [303] | 2021 | 30K frames | ✓ | onboard | 3D OD, OT | real |
| IPS3000+ [165] | 2021 | 14,198 frames | ✓ | V2X | 3D OD | real |
| AUTOMATUM [304] | 2021 | 30 hours | ✓ | drone | trajectory prediction | real |
| DurLAR [305] | 2021 | 100K frames | ✓ | onboard | depth estimation | real |
| Reasonable-Crowd [306] | 2021 | 92 scenarios | ✓ | onboard | driving behavior prediction | synthetic |
| MOTSynth [307] | 2021 | 768 driving sequences | ✓ | onboard | Pedestrian detection and tracking | synthetic |
| MAVD [308] | 2021 | 113,283 images | ✓ | onboard | OD and OT with sound | real |
| Multifog KITTI [309] | 2021 | 15K frames | ✗ | onboard | 3D OD | synthetic |
| Comap [310] | 2021 | 4,391 frames | ✓ | V2X | 3D OD | synthetic |
| R3 [311] | 2021 | 369 scenes | ✗ | onboard | out-of-distribution detection | real |
| WIBAM [312] | 2021 | 33,092 frames | ✓ | V2X | 3D OD | real |
| CeyMo [313] | 2021 | 2,887 frames | ✗ | onboard | road marking detection | real |
| RaidaR [314] | 2021 | 58,542 rainy street scenes | ✗ | onboard | SS | real |
| Fishyscapes [315] | 2021 | 1,030 frames | ✗ | onboard | SS, anomaly detection | real |
| RadarScenes [316] | 2021 | 40K radar frames | ✓ | onboard | OD, classification | real |
| ROAD [317] | 2021 | 122K frames | ✓ | onboard | OD, SS | real |
| All-in-One Drive [46] | 2021 | 100K | ✓ | onboard | (3D) OD, (3D) SS, trajectory prediction | real |
| PandaSet [166] | 2021 | 8,240 frames | ✓ | onboard | (3D) OD, SS, OT | real |
| SODA10M [318] | 2021 | 20K labeled images | ✓ | onboard | OD | real |
| PixSet [319] | 2021 | 29K point cloud frames | ✗ | onboard | (3D) OD | real |
| RoadObstacle21 [320] | 2021 | 327 scenes | ✗ | onboard | anomaly segmentation | synthetic |
| VIL-100 [321] | 2021 | 10K frames | ✗ | onboard | lane detection | real |
| OpenMPD [322] | 2021 | 15K frames | ✗ | onboard | (3D) OD, 3D OT, semantic segmentation | real |
| WADS [323] | 2021 | 1K point cloud frames | ✓ | onboard | SS | real |
| CCTSDB 2021 [324] | 2021 | 16,356 frames | ✗ | onboard | traffic sign detection | real |
| SemanticUSL [325] | 2021 | 1.2K frames | ✗ | onboard | domain adaptation 3D SS | real |
| CARLANE [326] | 2022 | 118K frames | ✓ | onboard | lane detection | synthetic |
| CrashD [327] | 2022 | 15,340 scenes | ✗ | onboard | 3D OD | synthetic |
| CODD [159] | 2022 | 108 sequences | ✓ | V2X | multi-agent SLAM | synthetic |
| CarlaScenes [328] | 2022 | 7 sequences | ✓ | onboard | (3D) SS, SLAM, depth estimation | synthetic |
| OPV2V [329] | 2022 | 11,464 frames | ✗ | V2X | onboard-to-onboard perception | synthetic |
| CARLA-WildLife [330] | 2022 | 26 videos | ✓ | onboard | out-of-distribution tracking | synthetic |
| SOS [330] | 2022 | 20 videos | ✓ | onboard | out-of-distribution tracking | real |
| RoadSaW [331] | 2022 | 720K frames | ✓ | onboard | Road surface and wetness estimation | real |
| I see you [332] | 2022 | 170 sequences, 340 trajectories | ✓ | V2X | OD | real |
| ASAP [333] | 2022 | 1.2M images | ✓ | onboard | online 3D OD | real |
| Amodal Cityscapes [334] | 2022 | 5K frames | ✗ | onboard | amodal SS | real |
| SynWoodScape [335] | 2022 | 80K frames | ✗ | onboard | (3D) OD, segmentation | synthetic |
| TJ4RadSet [336] | 2022 | 7,757 frames | ✓ | onboard | OD, OT | real |
| CODA [337] | 2022 | 1,500 frames | ✗ | onboard | corner case detection | real |
| LiDAR Snowfall [338] | 2022 | 7,385 point cloud frames | ✓ | onboard | 3D OD | synthetic |
| MUAD [339] | 2022 | 10.4K frames | ✓ | onboard | OD. SS, depth estimation | synthetic |
| AUTOCASTSIM [340] | 2022 | 52K frames | ✓ | V2X | (3D) OD, OT, SS | real |
| CARTI [341] | 2022 | 11K frames | ✓ | V2X | cooperative perception | synthetic |
| K-Lane [342] | 2022 | 15,382 frames | ✗ | onboard | lane detection | real |
| Ithaca365 [343] | 2022 | 7K frames | ✗ | onboard | 3D OD, SS, depth estimation | real |
| GLARE [344] | 2022 | 2,157 frames | ✗ | onboard | traffic sign detection | real |
| SUPS [345] | 2023 | 5K frames | ✓ | onboard | SS, depth estimation, SLAM | synthetic |
| Boreas [346] | 2023 | 7,111 frames | ✓ | onboard | (3D) OD, localization | real |
| Robo3D [347] | 2023 | 476K frames | ✓ | onboard | (3D) OD, 3D SS | real |
| ZOD [31] | 2023 | 100K frames | ✓ | onboard | (3D) OD, segmentation | real |
| K-Radar [348] | 2023 | 35K radar frames | ✗ | onboard | 3D OD, OT | real |
| aiMotive [349] | 2023 | 26,583 frames | ✓ | onboard | 3D OD, MOT | real |
| UrbanLaneGraph [350] | 2023 | around 5,220 km lane spans | ✓ | drone | lane graph estimation | real |
| WEDGE [351] | 2023 | 3,360 frames | ✗ | others | OD, classification | synthetic |
| OpenLane-V2 [352] | 2023 | 466K images | ✓ | onboard | lane detection, scene understanding | real |
| V2X-Seq (perception) [56] | 2023 | 15K frames | ✓ | V2X | cooperative perception | real |
| SSCBENCH [353] | 2023 | 66,913 frames | ✗ | onboard | semantic scene completion | real |
| RoadSC [354] | 2023 | 90,759 images | ✗ | onboard | road snow coverage estimation | real |
| V2X-Real [355] | 2024 | 171K images, 33K LiDAR point clouds | - | V2X | 3D OD | real |
| RCooper [356] | 2024 | 50K images, 30K LiDAR point clouds | ✓ | V2X | 3D OD, OT | real |
| FLIR [48] | - | 26,442 thermal frames | ✓ | onboard | thermal image OD | real |

## REFERENCES

[1] J. Mao, S. Shi, X. Wang, and H. Li, "3d object detection for autonomous driving: A comprehensive survey," *International Journal of Computer Vision*, pp. 1–55, 2023.

[2] G. Ros, L. Sellart, J. Materzynska, D. Vazquez, and A. M. Lopez, "The synthia dataset: A large collection of synthetic images for semantic segmentation of urban scenes," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3234–3243, 2016.

[3] A. Gaidon, Q. Wang, Y. Cabon, and E. Vig, "Virtual worlds as proxy for multi-object tracking analysis," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4340–4349, 2016.

[4] S. R. Richter, V. Vineet, S. Roth, and V. Koltun, "Playing for data: Ground truth from computer games," in *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part II 14*, pp. 102–118, Springer, 2016.

[5] S. R. Richter, Z. Hayder, and V. Koltun, "Playing for benchmarks," in *Proceedings of the IEEE International Conference on Computer Vision*, pp. 2213–2222, 2017.

[6] 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.

[7] D. Hendrycks, S. Basart, M. Mazeika, A. Zou, J. Kwon, M. Mostajabi, J. Steinhardt, and D. Song, "Scaling out-of-distribution detection for real-world settings," *arXiv preprint arXiv:1911.11132*, 2019.

[8] B. Hurl, K. Czarnecki, and S. Waslander, "Precise synthetic image and lidar (presil) dataset for autonomous vehicle perception," in *2019 IEEE Intelligent Vehicles Symposium (IV)*, pp. 2522–2529, IEEE, 2019.

[9] Y. Cabon, N. Murray, and M. Humenberger, "Virtual kitti 2," *arXiv preprint arXiv:2001.10773*, 2020.

[10] T. Sun, M. Segu, J. Postels, Y. Wang, L. Van Gool, B. Schiele, F. Tombari, and F. Yu, "Shift: a synthetic driving dataset for continuous multi-task domain adaptation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 21371–21382, 2022.

[11] R. Xu, H. Xiang, Z. Tu, X. Xia, M.-H. Yang, and J. Ma, "V2x-vit: Vehicle-to-everything cooperative perception with vision transformer," in *European conference on computer vision*, pp. 107–124, Springer, 2022.

[12] T. Wang, S. Kim, W. Ji, E. Xie, C. Ge, J. Chen, Z. Li, and P. Luo, "Deepaccident: A motion and accident prediction benchmark for v2x autonomous driving," *arXiv preprint arXiv:2304.01168*, 2023.

[13] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," in *2012 IEEE conference on computer vision and pattern recognition*, pp. 3354–3361, IEEE, 2012.

[14] 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.

[15] Z. Zhu, D. Liang, S. Zhang, X. Huang, B. Li, and S. Hu, "Traffic-sign detection and classification in the wild," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2110–2118, 2016.

[16] A. Rasouli, I. Kotseruba, and J. K. Tsotsos, "Are they going to cross? a benchmark dataset and baseline for pedestrian crosswalk behavior," in *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pp. 206–213, 2017.

[17] X. Huang, X. Cheng, Q. Geng, B. Cao, D. Zhou, P. Wang, Y. Lin, and R. Yang, "The apolloscape dataset for autonomous driving," in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pp. 954–960, 2018.

[18] R. Krajewski, J. Bock, L. Kloeker, and L. Eckstein, "The highd dataset: A drone dataset of naturalistic vehicle trajectories on german highways for validation of highly automated driving systems," in *2018 21st international conference on intelligent transportation systems (ITSC)*, pp. 2118–2125, IEEE, 2018.

[19] M.-F. Chang, J. Lambert, P. Sangkloy, J. Singh, S. Bak, A. Hartnett, D. Wang, P. Carr, S. Lucey, D. Ramanan, *et al.*, "Argoverse: 3d tracking and forecasting with rich maps," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 8748–8757, 2019.

[20] W. Zhan, L. Sun, D. Wang, H. Shi, A. Clausse, M. Naumann, J. Kummerle, H. Konigshof, C. Stiller, A. de La Fortelle, *et al.*, "Interaction dataset: An international, adversarial and cooperative motion dataset in interactive driving scenarios with semantic maps," *arXiv preprint arXiv:1910.03088*, 2019.

[21] A. Rasouli, I. Kotseruba, T. Kunic, and J. K. Tsotsos, "Pie: A large-scale dataset and models for pedestrian intention estimation and trajectory prediction," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 6262–6271, 2019.

[22] 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.

[23] 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.

[24] 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.

[25] J. Geyer, Y. Kassahun, M. Mahmudi, X. Ricou, R. Durgesh, A. S. Chung, L. Hauswald, V. H. Pham, M. Mühlegg, S. Dorn, *et al.*, "A2d2: Audi autonomous driving dataset," *arXiv preprint arXiv:2004.06320*, 2020.

[26] M. Sheeny, E. De Pellegrin, S. Mukherjee, A. Ahrabian, S. Wang, and A. Wallace, "Radiate: A radar dataset for automotive perception in bad weather," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 1–7, IEEE, 2021.

[27] H. Yu, Y. Luo, M. Shu, Y. Huo, Z. Yang, Y. Shi, Z. Guo, H. Li, X. Hu, J. Yuan, *et al.*, "Dair-v2x: A large-scale dataset for vehicle-infrastructure cooperative 3d object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 21361–21370, 2022.

[28] B. Wilson, W. Qi, T. Agarwal, J. Lambert, J. Singh, S. Khandelwal, B. Pan, R. Kumar, A. Hartnett, J. K. Pontes, *et al.*, "Argoverse 2: Next generation datasets for self-driving perception and forecasting," *arXiv preprint arXiv:2301.00493*, 2023.

[29] Y. Li, Z. Li, S. Teng, Y. Zhang, Y. Zhou, Y. Zhu, D. Cao, B. Tian, Y. Ai, Z. Xuanyuan, *et al.*, "Automine: An unmanned mine dataset," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 21308–21317, 2022.

[30] J. Mao, M. Niu, C. Jiang, H. Liang, J. Chen, X. Liang, Y. Li, C. Ye, W. Zhang, Z. Li, *et al.*, "One million scenes for autonomous driving: Once dataset," *arXiv preprint arXiv:2106.11037*, 2021.

[31] M. Alibeigi, W. Ljungbergh, A. Tonderski, G. Hess, A. Lilja, C. Lindström, D. Motorniuk, J. Fu, J. Widahl, and C. Petersson, "Zenseact open dataset: A large-scale and diverse multimodal dataset for autonomous driving," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 20178–20188, 2023.

[32] H. Yin and C. Berger, "When to use what data set for your self-driving car algorithm: An overview of publicly available driving datasets," in *2017 IEEE 20th International Conference on Intelligent Transportation Systems (ITSC)*, pp. 1–8, IEEE, 2017.

[33] Y. Kang, H. Yin, and C. Berger, "Test your self-driving algorithm: An overview of publicly available driving datasets and virtual testing environments," *IEEE Transactions on Intelligent Vehicles*, vol. 4, no. 2, pp. 171–185, 2019.

[34] J. Guo, U. Kurup, and M. Shah, "Is it safe to drive? an overview of factors, metrics, and datasets for driveability assessment in autonomous driving," *IEEE Transactions on Intelligent Transportation Systems*, vol. 21, no. 8, pp. 3135–3151, 2019.

[35] J. Janai, F. Güney, A. Behl, A. Geiger, *et al.*, "Computer vision for autonomous vehicles: Problems, datasets and state of the art," *Foundations and Trends® in Computer Graphics and Vision*, vol. 12, no. 1–3, pp. 1–308, 2020.

[36] W. Liu, Q. Dong, P. Wang, G. Yang, L. Meng, Y. Song, Y. Shi, and Y. Xue, "A survey on autonomous driving datasets," in *2021 8th International Conference on Dependable Systems and Their Applications (DSA)*, pp. 399–407, IEEE, 2021.

[37] H. Li, Y. Li, H. Wang, J. Zeng, P. Cai, H. Xu, D. Lin, J. Yan, F. Xu, L. Xiong, *et al.*, "Open-sourced data ecosystem in autonomous driving: the present and future," *arXiv preprint arXiv:2312.03408*, 2023.

[38] D. Bogdoll, F. Schreyer, and J. M. Zöllner, "Ad-datasets: a meta-collection of data sets for autonomous driving," *arXiv preprint arXiv:2202.01909*, 2022.

[39] D. Bogdoll, S. Uhlemeyer, K. Kowol, and J. M. Zöllner, "Perception datasets for anomaly detection in autonomous driving: A survey," *arXiv preprint arXiv:2302.02790*, 2023.

[40] Z. Song, Z. He, X. Li, Q. Ma, R. Ming, Z. Mao, H. Pei, L. Peng, J. Hu, D. Yao, *et al.*, "Synthetic datasets for autonomous driving: A survey," *arXiv preprint arXiv:2304.12205*, 2023.

[41] B. Gao, Y. Pan, C. Li, S. Geng, and H. Zhao, "Are we hungry for 3d lidar data for semantic segmentation? a survey of datasets and methods," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 7, pp. 6063–6081, 2021.

[42] Y. Wang, Z. Han, Y. Xing, S. Xu, and J. Wang, "A survey on datasets for decision-making of autonomous vehicle," *arXiv preprint arXiv:2306.16784*, 2023.

[43] S. Teng, X. Hu, P. Deng, B. Li, Y. Li, Y. Ai, D. Yang, L. Li, Z. Xuanyuan, F. Zhu, *et al.*, "Motion planning for autonomous driving: The state of the art and future perspectives," *IEEE Transactions on Intelligent Vehicles*, 2023.

[44] L. Chen, Y. Li, C. Huang, B. Li, Y. Xing, D. Tian, L. Li, Z. Hu, X. Na, Z. Li, *et al.*, "Milestones in autonomous driving and intelligent vehicles: Survey of surveys," *IEEE Transactions on Intelligent Vehicles*, vol. 8, no. 2, pp. 1046–1056, 2022.

[45] B. Kitchenham, "Procedures for performing systematic reviews," *Keele, UK, Keele University*, vol. 33, no. 2004, pp. 1–26, 2004.

[46] X. Weng, Y. Man, J. Park, Y. Yuan, M. O'Toole, and K. M. Kitani, "All-in-one drive: A comprehensive perception dataset with high-density long-range point clouds," 2023.

[47] M. Gehrig, W. Aarents, D. Gehrig, and D. Scaramuzza, "Dsec: A stereo event camera dataset for driving scenarios," *IEEE Robotics and Automation Letters*, vol. 6, no. 3, pp. 4947–4954, 2021.

[48] "Flir: https://www.flir.com/oem/adas/adas-dataset-form/."

[49] S. Grigorescu, B. Trasnea, T. Cocias, and G. Macesanu, "A survey of deep learning techniques for autonomous driving," *Journal of Field Robotics*, vol. 37, no. 3, pp. 362–386, 2020.

[50] Y. Li and J. Ibanez-Guzman, "Lidar for autonomous driving: The principles, challenges, and trends for automotive lidar and perception systems," *IEEE Signal Processing Magazine*, vol. 37, no. 4, pp. 50–61, 2020.

[51] T. Zhou, M. Yang, K. Jiang, H. Wong, and D. Yang, "Mmw radar-based technologies in autonomous driving: A review," *Sensors*, vol. 20, no. 24, p. 7283, 2020.

[52] G. Chen, H. Cao, J. Conradt, H. Tang, F. Rohrbein, and A. Knoll, "Event-based neuromorphic vision for autonomous driving: A paradigm shift for bio-inspired visual sensing and perception," *IEEE Signal Processing Magazine*, vol. 37, no. 4, pp. 34–49, 2020.

[53] R. Gade and T. B. Moeslund, "Thermal cameras and applications: a survey," *Machine vision and applications*, vol. 25, pp. 245–262, 2014.

[54] E. Yurtsever, J. Lambert, A. Carballo, and K. Takeda, "A survey of autonomous driving: Common practices and emerging technologies," *IEEE access*, vol. 8, pp. 58443–58469, 2020.

[55] T. Huang, J. Liu, X. Zhou, D. C. Nguyen, M. R. Azghadi, Y. Xia, Q.-L. Han, and S. Sun, "V2x cooperative perception for autonomous driving: Recent advances and challenges," *arXiv preprint arXiv:2310.03525*, 2023.

[56] H. Yu, W. Yang, H. Ruan, Z. Yang, Y. Tang, X. Gao, X. Hao, Y. Shi, Y. Pan, N. Sun, *et al.*, "V2x-seq: A large-scale sequential dataset for vehicle-infrastructure cooperative perception and forecasting," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5486–5495, 2023.

[57] 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.

[58] Y. Liao, J. Xie, and A. Geiger, "Kitti-360: A novel dataset and benchmarks for urban scene understanding in 2d and 3d," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 3, pp. 3292–3310, 2022.

[59] A. Palazzi, D. Abati, F. Solera, R. Cucchiara, *et al.*, "Predicting the driver's focus of attention: the dr (eye) ve project," *IEEE transactions on pattern analysis and machine intelligence*, vol. 41, no. 7, pp. 1720–1733, 2018.

[60] W. Zimmer, G. A. Wardana, S. Sritharan, X. Zhou, R. Song, and A. C. Knoll, "Tumtraf v2x cooperative perception dataset," in *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, IEEE, 2024.

[61] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pp. 740–755, Springer, 2014.

[62] A. Garcia-Garcia, S. Orts-Escolano, S. Oprea, V. Villena-Martinez, P. Martinez-Gonzalez, and J. Garcia-Rodriguez, "A survey on deep learning techniques for image and video semantic segmentation," *Applied Soft Computing*, vol. 70, pp. 41–65, 2018.

[63] W. Luo, J. Xing, A. Milan, X. Zhang, W. Liu, and T.-K. Kim, "Multiple object tracking: A literature review," *Artificial intelligence*, vol. 293, p. 103448, 2021.

[64] S. Guo, S. Wang, Z. Yang, L. Wang, H. Zhang, P. Guo, Y. Gao, and J. Guo, "A review of deep learning-based visual multi-object tracking algorithms for autonomous driving," *Applied Sciences*, vol. 12, no. 21, p. 10741, 2022.

[65] Q. Li, Y. Wang, Y. Wang, and H. Zhao, "Hdmapnet: An online hd map construction and evaluation framework," in *2022 International Conference on Robotics and Automation (ICRA)*, pp. 4628–4634, IEEE, 2022.

[66] J. Lambert and J. Hays, "Trust, but verify: Cross-modality fusion for hd map change detection," *arXiv preprint arXiv:2212.07312*, 2022.

[67] J. Sturm, N. Engelhard, F. Endres, W. Burgard, and D. Cremers, "A benchmark for the evaluation of rgb-d slam systems," in *2012 IEEE/RSJ international conference on intelligent robots and systems*, pp. 573–580, IEEE, 2012.

[68] Y. Huang, J. Du, Z. Yang, Z. Zhou, L. Zhang, and H. Chen, "A survey on trajectory-prediction methods for autonomous driving," *IEEE Transactions on Intelligent Vehicles*, vol. 7, no. 3, pp. 652–674, 2022.

[69] S. Mozaffari, O. Y. Al-Jarrah, M. Dianati, P. Jennings, and A. Mouzakitis, "Deep learning-based vehicle behavior prediction for autonomous driving applications: A review," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 1, pp. 33–47, 2020.

[70] W. Ding, J. Chen, and S. Shen, "Predicting vehicle behaviors over an extended horizon using behavior interaction network," in *2019 international conference on robotics and automation (ICRA)*, pp. 8634–8640, IEEE, 2019.

[71] N. Sharma, C. Dhiman, and S. Indu, "Pedestrian intention prediction for autonomous vehicles: A comprehensive survey," *Neurocomputing*, 2022.

[72] B. Paden, M. Čáp, S. Z. Yong, D. Yershov, and E. Frazzoli, "A survey of motion planning and control techniques for self-driving urban vehicles," *IEEE Transactions on intelligent vehicles*, vol. 1, no. 1, pp. 33–55, 2016.

[73] A. Tampuu, T. Matiisen, M. Semikin, D. Fishman, and N. Muhammad, "A survey of end-to-end driving: Architectures and training methods," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 33, no. 4, pp. 1364–1384, 2020.

[74] 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.

[75] X. Pan, J. Shi, P. Luo, X. Wang, and X. Tang, "Spatial as deep: Spatial cnn for traffic scene understanding," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, 2018.

[76] 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.

[77] C. Sakaridis, D. Dai, and L. Van Gool, "Semantic foggy scene understanding with synthetic data," *International Journal of Computer Vision*, vol. 126, pp. 973–992, 2018.

[78] G. J. Brostow, J. Fauqueur, and R. Cipolla, "Semantic object classes in video: A high-definition ground truth database," *Pattern Recognition Letters*, vol. 30, no. 2, pp. 88–97, 2009.

[79] G. Varma, A. Subramanian, A. Namboodiri, M. Chandraker, and C. Jawahar, "Idd: A dataset for exploring problems of autonomous navigation in unconstrained environments," in *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pp. 1743–1751, IEEE, 2019.

[80] C. Sakaridis, D. Dai, and L. Van Gool, "Acdc: The adverse conditions dataset with correspondences for semantic driving scene understanding," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 10765–10775, 2021.

[81] S. Houben, J. Stallkamp, J. Salmen, M. Schlipsing, and C. Igel, "Detection of traffic signs in real-world images: The german traffic sign detection benchmark," in *The 2013 international joint conference on neural networks (IJCNN)*, pp. 1–8, Ieee, 2013.

[82] 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.

[83] M. Bijelic, T. Gruber, F. Mannan, F. Kraus, W. Ritter, K. Dietmayer, and F. Heide, "Seeing through fog without seeing fog: Deep multimodal sensor fusion in unseen adverse weather," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11682–11692, 2020.

[84] C. Creß, W. Zimmer, L. Strand, M. Fortkord, S. Dai, V. Lakshminarasimhan, and A. Knoll, "A9-dataset: Multi-sensor infrastructure-based dataset for mobility research," in *2022 IEEE Intelligent Vehicles Symposium (IV)*, pp. 965–970, IEEE, 2022.

[85] W. Zimmer, C. Creß, H. T. Nguyen, and A. C. Knoll, "Tumtraf intersection dataset: All you need for urban 3d camera-lidar roadside perception," in *2023 IEEE 26th International Conference on Intelligent Transportation Systems (ITSC)*, pp. 1030–1037, IEEE, 2023.

[86] C. Creß, W. Zimmer, N. Purschke, B. N. Doan, V. Lakshminarasimhan, L. Strand, and A. C. Knoll, "Tumtraf event: Calibration and fusion resulting in a dataset for roadside event-based and rgb cameras," *arXiv preprint arXiv:2401.08474*, 2024.

[87] R. Xu, X. Xia, J. Li, H. Li, S. Zhang, Z. Tu, Z. Meng, H. Xiang, X. Dong, R. Song, *et al.*, "V2v4real: A real-world large-scale dataset for vehicle-to-vehicle cooperative perception," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 13712–13722, 2023.

[88] X. Ye, M. Shu, H. Li, Y. Shi, Y. Li, G. Wang, X. Tan, and E. Ding, "Rope3d: The roadside perception dataset for autonomous driving and monocular 3d object detection task," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 21341–21350, 2022.

[89] Y. Li, D. Ma, Z. An, Z. Wang, Y. Zhong, S. Chen, and C. Feng, "V2x-sim: Multi-agent collaborative perception dataset and benchmark for autonomous driving," *IEEE Robotics and Automation Letters*, vol. 7, no. 4, pp. 10914–10921, 2022.

[90] T.-H. Wang, S. Manivasagam, M. Liang, B. Yang, W. Zeng, and R. Urtasun, "V2vnet: Vehicle-to-vehicle communication for joint perception and prediction," in *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*, pp. 605–621, Springer, 2020.

[91] Q. Chen, S. Tang, Q. Yang, and S. Fu, "Cooper: Cooperative perception for connected autonomous vehicles based on 3d point clouds," in *2019 IEEE 39th International Conference on Distributed Computing Systems (ICDCS)*, pp. 514–524, IEEE, 2019.

[92] E. Arnold, M. Dianati, R. de Temple, and S. Fallah, "Cooperative perception for 3d object detection in driving scenarios using infrastructure sensors," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 3, pp. 1852–1864, 2020.

[93] S. Busch, C. Koetsier, J. Axmann, and C. Brenner, "Lumpi: The leibniz university multi-perspective intersection dataset," in *2022 IEEE Intelligent Vehicles Symposium (IV)*, pp. 1127–1134, IEEE, 2022.

[94] D. Du, Y. Qi, H. Yu, Y. Yang, K. Duan, G. Li, W. Zhang, Q. Huang, and Q. Tian, "The unmanned aerial vehicle benchmark: Object detection and tracking," in *Proceedings of the European conference on computer vision (ECCV)*, pp. 370–386, 2018.

[95] Y. Sun, B. Cao, P. Zhu, and Q. Hu, "Drone-based rgb-infrared cross-modality vehicle detection via uncertainty-aware learning," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 32, no. 10, pp. 6700–6713, 2022.

[96] G. Neuhold, T. Ollmann, S. Rota Bulo, and P. Kontschieder, "The mapillary vistas dataset for semantic understanding of street scenes," in *Proceedings of the IEEE international conference on computer vision*, pp. 4990–4999, 2017.

[97] Y. Xiang, R. Mottaghi, and S. Savarese, "Beyond pascal: A benchmark for 3d object detection in the wild," in *IEEE winter conference on applications of computer vision*, pp. 75–82, IEEE, 2014.

[98] O. Zendel, K. Honauer, M. Murschitz, D. Steininger, and G. F. Dominguez, "Wilddash-creating hazard-aware benchmarks," in *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 402–416, 2018.

[99] S. Wang, M. Bai, G. Mattyus, H. Chu, W. Luo, B. Yang, J. Liang, J. Cheverie, S. Fidler, and R. Urtasun, "Torontocity: Seeing the world with a million eyes," *arXiv preprint arXiv:1612.00423*, 2016.

[100] M. A. Kenk and M. Hassaballah, "Dawn: vehicle detection in adverse weather nature dataset," *arXiv preprint arXiv:2008.05402*, 2020.

[101] K. Lis, K. Nakka, P. Fua, and M. Salzmann, "Detecting the unexpected via image resynthesis," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 2152–2161, 2019.

[102] P. Cong, X. Zhu, F. Qiao, Y. Ren, X. Peng, Y. Hou, L. Xu, R. Yang, D. Manocha, and Y. Ma, "Stcrowd: A multimodal dataset for pedestrian perception in crowded scenes. 2022 ieee," in *CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 19576–19585, 2022.

[103] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes (voc) challenge," *International journal of computer vision*, vol. 88, pp. 303–338, 2010.

[104] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255, Ieee, 2009.

[105] A. Jain, H. S. Koppula, S. Soh, B. Raghavan, A. Singh, and A. Saxena, "Brain4cars: Car that knows before you do via sensory-fusion deep learning architecture," *arXiv preprint arXiv:1601.00740*, 2016.

[106] A. Zyner, S. Worrall, and E. Nebot, "Naturalistic driver intention and path prediction using recurrent neural networks," *IEEE transactions on intelligent transportation systems*, vol. 21, no. 4, pp. 1584–1594, 2019.

[107] M. Martin, A. Roitberg, M. Haurilet, M. Horne, S. Reiß, M. Voit, and R. Stiefelhagen, "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*, pp. 2801–2810, 2019.

[108] Y. Chen, J. Wang, J. Li, C. Lu, Z. Luo, H. Xue, and C. Wang, "Dbnet: A large-scale dataset for driving behavior learning," *Retrieved August*, vol. 9, p. 2020, 2019.

[109] B. Toghi, D. Grover, M. Razzaghpour, R. Jain, R. Valiente, M. Zaman, G. Shah, and Y. P. Fallah, "A maneuver-based urban driving dataset and model for cooperative vehicle applications," in *2020 IEEE 3rd Connected and Automated Vehicles Symposium (CAVS)*, pp. 1–6, IEEE, 2020.

[110] J. Bock, R. Krajewski, T. Moers, S. Runde, L. Vater, and L. Eckstein, "The ind dataset: A drone dataset of naturalistic road user trajectories at german intersections," in *2020 IEEE Intelligent Vehicles Symposium (IV)*, pp. 1929–1934, IEEE, 2020.

[111] A. Rasouli, T. Yau, P. Lakner, S. Malekmohammadi, M. Rohani, and J. Luo, "Pepscenes: A novel dataset and baseline for pedestrian action prediction in 3d," *arXiv preprint arXiv:2012.07773*, 2020.

[112] A. Breuer, J.-A. Termöhlen, S. Homoceanu, and T. Fingscheidt, "opendd: A large-scale roundabout drone dataset," in *2020 IEEE 23rd International Conference on Intelligent Transportation Systems (ITSC)*, pp. 1–6, IEEE, 2020.

[113] H. Caesar, J. Kabzan, K. S. Tan, W. K. Fong, E. Wolff, A. Lang, L. Fletcher, O. Beijbom, and S. Omari, "nuplan: A closed-loop ml-based planning benchmark for autonomous vehicles," *arXiv preprint arXiv:2106.11810*, 2021.

[114] R. Guesdon, C. Crispim-Junior, and L. Tougne, "Dripe: A dataset for human pose estimation in real-world driving settings," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 2865–2874, 2021.

[115] S. Ghosh, A. Dhall, G. Sharma, S. Gupta, and N. Sebe, "Speak2label: Using domain knowledge for creating a large scale driver gaze zone estimation dataset," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 2896–2905, 2021.

[116] Y. Shen, N. Wijayaratne, P. Sriram, A. Hasan, P. Du, and K. Driggs-Campbell, "Cocatt: A cognitive-conditioned driver attention dataset," in *2022 IEEE 25th International Conference on Intelligent Transportation Systems (ITSC)*, pp. 32–39, IEEE, 2022.

[117] T. Moers, L. Vater, R. Krajewski, J. Bock, A. Zlocki, and L. Eckstein, "The exid dataset: A real-world trajectory dataset of highly interactive highway scenarios in germany," in *2022 IEEE Intelligent Vehicles Symposium (IV)*, pp. 958–964, IEEE, 2022.

[118] L. Gressenbuch, K. Esterle, T. Kessler, and M. Althoff, "Mona: The munich motion dataset of natural driving," in *2022 IEEE 25th International Conference on Intelligent Transportation Systems (ITSC)*, pp. 2093–2100, IEEE, 2022.

[119] X. Tian, T. Jiang, L. Yun, Y. Mao, H. Yang, Y. Wang, Y. Wang, and H. Zhao, "Occ3d: A large-scale 3d occupancy prediction benchmark for autonomous driving," *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[120] V. Ramanishka, Y.-T. Chen, T. Misu, and K. Saenko, "Toward driving scene understanding: A dataset for learning driver behavior and causal reasoning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7699–7707, 2018.

[121] J. Xue, J. Fang, T. Li, B. Zhang, P. Zhang, Z. Ye, and J. Dou, "Blvd: Building a large-scale 5d semantics benchmark for autonomous driving," in *2019 International Conference on Robotics and Automation (ICRA)*, pp. 6685–6691, IEEE, 2019.

[122] R. Krajewski, T. Moers, J. Bock, L. Vater, and L. Eckstein, "The round dataset: A drone dataset of road user trajectories at roundabouts in

germany," in *2020 IEEE 23rd International Conference on Intelligent Transportation Systems (ITSC)*, pp. 1–6, IEEE, 2020.

[123] R. Izquierdo, A. Quintanar, I. Parra, D. Fernández-Llorca, and M. Sotelo, "The prevention dataset: a novel benchmark for prediction of vehicles intentions," in *2019 IEEE Intelligent Transportation Systems Conference (ITSC)*, pp. 3114–3121, IEEE, 2019.

[124] G. Yang, X. Song, C. Huang, Z. Deng, J. Shi, and B. Zhou, "Driving-stereo: A large-scale dataset for stereo matching in autonomous driving scenarios," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 899–908, 2019.

[125] J. Houston, G. Zuidhof, L. Bergamini, Y. Ye, L. Chen, A. Jain, S. Omari, V. Iglovikov, and P. Ondruska, "One thousand and one hours: Self-driving motion prediction dataset," in *Conference on Robot Learning*, pp. 409–418, PMLR, 2021.

[126] H. Girase, H. Gang, S. Malla, J. Li, A. Kanehara, K. Mangalam, and C. Choi, "Loki: Long term and key intentions for trajectory prediction," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 9803–9812, 2021.

[127] A. Prabu, N. Ranjan, L. Li, R. Tian, S. Chien, Y. Chen, and R. Sherony, "Scendd: A scenario-based naturalistic driving dataset," in *2022 IEEE 25th International Conference on Intelligent Transportation Systems (ITSC)*, pp. 4363–4368, IEEE, 2022.

[128] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun, "Carla: An open urban driving simulator," in *Conference on robot learning*, pp. 1–16, PMLR, 2017.

[129] J. Binas, D. Neil, S.-C. Liu, and T. Delbruck, "Ddd17: End-to-end davis driving dataset," *arXiv preprint arXiv:1711.01458*, 2017.

[130] B. C. Russell, A. Torralba, K. P. Murphy, and W. T. Freeman, "Labelme: a database and web-based tool for image annotation," *International journal of computer vision*, vol. 77, pp. 157–173, 2008.

[131] B.-L. Wang, C.-T. King, and H.-K. Chu, "A semi-automatic video labeling tool for autonomous driving based on multi-object detector and tracker," in *2018 sixth international symposium on computing and networking (CANDAR)*, pp. 201–206, IEEE, 2018.

[132] "https://www.cs.columbia.edu/ vondrick/vatic/,"

[133] N. Manikandan and K. Ganesan, "Deep learning based automatic video annotation tool for self-driving car," *arXiv preprint arXiv:1904.12618*, 2019.

[134] D. Schörkhuber, F. Groh, and M. Gelautz, "Bounding box propagation for semi-automatic video annotation of nighttime driving scenes," in *2021 12th International Symposium on Image and Signal Processing and Analysis (ISPA)*, pp. 131–137, IEEE, 2021.

[135] Q. Meng, W. Wang, T. Zhou, J. Shen, Y. Jia, and L. Van Gool, "Towards a weakly supervised framework for 3d point cloud object detection and annotation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 8, pp. 4454–4468, 2021.

[136] D. Zhang, D. Liang, Z. Zou, J. Li, X. Ye, Z. Liu, X. Tan, and X. Bai, "A simple vision transformer for weakly semi-supervised 3d object detection," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 8373–8383, 2023.

[137] J. Xie, M. Kiefel, M.-T. Sun, and A. Geiger, "Semantic instance annotation of street scenes by 3d to 2d label transfer," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pp. 3688–3697, 2016.

[138] W. Zimmer, A. Rangesh, and M. Trivedi, "3d bat: A semi-automatic, web-based 3d annotation toolbox for full-surround, multi-modal data streams," in *2019 IEEE Intelligent Vehicles Symposium (IV)*, pp. 1816–1821, IEEE, 2019.

[139] D. Barnes, W. Maddern, and I. Posner, "Find your own way: Weakly-supervised segmentation of path proposals for urban autonomy," in *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 203–210, IEEE, 2017.

[140] F. Saleh, M. S. Aliakbarian, M. Salzmann, L. Petersson, S. Gould, and J. M. Alvarez, "Built-in foreground/background prior for weakly-supervised semantic segmentation," in *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part VIII 14*, pp. 413–432, Springer, 2016.

[141] A. Petrovai, A. D. Costea, and S. Nedevschi, "Semi-automatic image annotation of street scenes," in *2017 IEEE intelligent vehicles symposium (IV)*, pp. 448–455, IEEE, 2017.

[142] D. Acuna, H. Ling, A. Kar, and S. Fidler, "Efficient interactive annotation of segmentation datasets with polygon-rnn++," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pp. 859–868, 2018.

[143] L. Castrejon, K. Kundu, R. Urtasun, and S. Fidler, "Annotating object instances with a polygon-rnn," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 5230–5238, 2017.

[144] Z. Chen, Q. Liao, Z. Wang, Y. Liu, and M. Liu, "Image detector based automatic 3d data labeling and training for vehicle detection on point cloud," in *2019 IEEE Intelligent Vehicles Symposium (IV)*, pp. 1408–1413, IEEE, 2019.

[145] H. Luo, C. Wang, C. Wen, Z. Chen, D. Zai, Y. Yu, and J. Li, "Semantic labeling of mobile lidar point clouds via active learning and higher order mrf," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 56, no. 7, pp. 3631–3644, 2018.

[146] M. Liu, Y. Zhou, C. R. Qi, B. Gong, H. Su, and D. Anguelov, "Less: Label-efficient semantic segmentation for lidar point clouds," in *European Conference on Computer Vision*, pp. 70–89, Springer, 2022.

[147] M. Wang, T. Ganjineh, and R. Rojas, "Action annotated trajectory generation for autonomous maneuvers on structured road networks," in *The 5th International conference on automation, robotics and applications*, pp. 67–72, IEEE, 2011.

[148] S. Moosavi, B. Omidvar-Tehrani, R. B. Craig, and R. Ramnath, "Annotation of car trajectories based on driving patterns," *arXiv preprint arXiv:1705.05219*, 2017.

[149] S. Jarl, L. Aronsson, S. Rahrovani, and M. H. Chehreghani, "Active learning of driving scenario trajectories," *Engineering Applications of Artificial Intelligence*, vol. 113, p. 104972, 2022.

[150] O. Styles, A. Ross, and V. Sanchez, "Forecasting pedestrian trajectory with machine-annotated training data," in *2019 IEEE Intelligent Vehicles Symposium (IV)*, pp. 716–721, IEEE, 2019.

[151] B. Wymann, E. Espié, C. Guionneau, C. Dimitrakakis, R. Coulom, and A. Sumner, "Torcs, the open racing car simulator," *Software available at http://torcs. sourceforge. net*, vol. 4, no. 6, p. 2, 2000.

[152] C. Chen, A. Seff, A. Kornhauser, and J. Xiao, "Deepdriving: Learning affordance for direct perception in autonomous driving," in *Proceedings of the IEEE international conference on computer vision*, pp. 2722–2730, 2015.

[153] D. Krajzewicz, J. Erdmann, M. Behrisch, and L. Bieker, "Recent development and applications of sumo-simulation of urban mobility," *International journal on advances in systems and measurements*, vol. 5, no. 3&4, 2012.

[154] S. Shah, D. Dey, C. Lovett, and A. Kapoor, "Airsim: High-fidelity visual and physical simulation for autonomous vehicles," in *Field and Service Robotics: Results of the 11th International Conference*, pp. 621–635, Springer, 2018.

[155] "https://www.nvidia.com/en-us/self-driving-cars/simulation/,"

[156] A. Juliani, V.-P. Berges, E. Teng, A. Cohen, J. Harper, C. Elion, C. Goy, Y. Gao, H. Henry, M. Mattar, *et al.*, "Unity: A general platform for intelligent agents," *arXiv preprint arXiv:1809.02627*, 2018.

[157] P. Krähenbühl, "Free supervision from video games," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2955–2964, 2018.

[158] R. Song, R. Xu, A. Festag, J. Ma, and A. Knoll, "Fedbevt: Federated learning bird's eye view perception transformer in road traffic systems," *IEEE Transactions on Intelligent Vehicles*, 2023.

[159] E. Arnold, S. Mozaffari, and M. Dianati, "Fast and robust registration of partially overlapping point clouds," *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 1502–1509, 2021.

[160] F. S. Saleh, M. S. Aliakbarian, M. Salzmann, L. Petersson, and J. M. Alvarez, "Effective use of synthetic data for urban scene semantic segmentation," in *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 84–100, 2018.

[161] X. Li, P. Ye, J. Li, Z. Liu, L. Cao, and F.-Y. Wang, "From features engineering to scenarios engineering for trustworthy ai: I&i, c&c, and v&v," *IEEE Intelligent Systems*, vol. 37, no. 4, pp. 18–26, 2022.

[162] R. M. Monarch, *Human-in-the-Loop Machine Learning: Active learning and annotation for human-centered AI*. Simon and Schuster, 2021.

[163] H.-M. Heyn, K. M. Habibullah, E. Knauss, J. Horkoff, M. Borg, A. Knauss, and P. J. Li, "Automotive perception software development: An empirical investigation into data, annotation, and ecosystem challenges," *arXiv preprint arXiv:2303.05947*, 2023.

[164] O. Inel and L. Aroyo, "Validation methodology for expert-annotated datasets: Event annotation case study," in *2nd Conference on Language, Data and Knowledge (LDK 2019)*, Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2019.

[165] H. Wang, X. Zhang, J. Li, Z. Li, L. Yang, S. Pan, and Y. Deng, "Ips300+: a challenging multimodal dataset for intersection perception system," *arXiv preprint arXiv:2106.02781*, 2021.

[166] P. Xiao, Z. Shao, S. Hao, Z. Zhang, X. Chai, J. Jiao, Z. Li, J. Wu, K. Sun, K. Jiang, *et al.*, "Pandaset: Advanced sensor suite dataset for autonomous driving," in *2021 IEEE International Intelligent Transportation Systems Conference (ITSC)*, pp. 3095–3101, IEEE, 2021.

[167] D. Barnes, M. Gadd, P. Murcutt, P. Newman, and I. Posner, "The oxford radar robotcar dataset: A radar extension to the oxford robotcar dataset," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 6433–6438, IEEE, 2020.

[168] M. Pitropov, D. E. Garcia, J. Rebello, M. Smart, C. Wang, K. Czarnecki, and S. Waslander, "Canadian adverse driving conditions dataset," *The International Journal of Robotics Research*, vol. 40, no. 4-5, pp. 681–690, 2021.

[169] E. Li, S. Wang, C. Li, D. Li, X. Wu, and Q. Hao, "Sustech points: A portable 3d point cloud interactive annotation platform system," in *2020 IEEE Intelligent Vehicles Symposium (IV)*, pp. 1108–1115, IEEE, 2020.

[170] Y. Chen, J. Liu, X. Zhang, X. Qi, and J. Jia, "Voxelnext: Fully sparse voxelnet for 3d object detection and tracking," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 21674–21683, 2023.

[171] Y. Li, Y. Chen, X. Qi, Z. Li, J. Sun, and J. Jia, "Unifying voxel-based representation with transformer for 3d object detection," *Advances in Neural Information Processing Systems*, vol. 35, pp. 18442–18455, 2022.

[172] X. Bai, Z. Hu, X. Zhu, Q. Huang, Y. Chen, H. Fu, and C.-L. Tai, "Transfusion: Robust lidar-camera fusion for 3d object detection with transformers," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 1090–1099, 2022.

[173] Y. Cui, W. Ren, X. Cao, and A. Knoll, "Focal network for image restoration," in *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 13001–13011, 2023.

[174] Y. Cui, W. Ren, X. Cao, and A. Knoll, "Image restoration via frequency selection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.

[175] 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.

[176] 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.

[177] 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.

[178] 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.

[179] Y. Xu, X. Yang, L. Gong, H.-C. Lin, T.-Y. Wu, Y. Li, and N. Vasconcelos, "Explainable object-induced action decision for autonomous vehicles," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9523–9532, 2020.

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

[181] 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.

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

[183] 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.

[184] 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.

[185] 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.

[186] 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.

[187] Y. Ma, C. Cui, X. Cao, W. Ye, P. Liu, J. Lu, A. Abdelraouf, R. Gupta, K. Han, A. Bera, *et al.*, "Lampilot: An open benchmark dataset for autonomous driving with language model programs," *arXiv preprint arXiv:2312.04372*, 2023.

[188] H. Shao, Y. Hu, L. Wang, S. L. Waslander, Y. Liu, and H. Li, "Lmdrive: Closed-loop end-to-end driving with large language models," *arXiv preprint arXiv:2312.07488*, 2023.

[189] Y. Inoue, Y. Yada, K. Tanahashi, and Y. Yamaguchi, "Nuscenes-mqa: Integrated evaluation of captions and qa for autonomous driving datasets using markup annotations," *arXiv preprint arXiv:2312.06352*, 2023.

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

[191] C. Sima, K. Renz, K. Chitta, L. Chen, H. Zhang, C. Xie, P. Luo, A. Geiger, and H. Li, "Drivelm: Driving with graph visual question answering," *arXiv preprint arXiv:2312.14150*, 2023.

[192] S. Yang, J. Liu, R. Zhang, M. Pan, Z. Guo, X. Li, Z. Chen, P. Gao, Y. Guo, and S. Zhang, "Lidar-llm: Exploring the potential of large language models for 3d lidar understanding," *arXiv preprint arXiv:2312.14074*, 2023.

[193] 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.

[194] Y. Inoue, Y. Yada, K. Tanahashi, and Y. Yamaguchi, "Nuscenes-mqa: Integrated evaluation of captions and qa for autonomous driving datasets using markup annotations," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 930–938, 2024.

[195] S. Park, M. Lee, J. Kang, H. Choi, Y. Park, J. Cho, A. Lee, and D. Kim, "Vlaad: Vision and language assistant for autonomous driving," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 980–987, 2024.

[196] Y. Hu, J. Yang, L. Chen, K. Li, C. Sima, X. Zhu, S. Chai, S. Du, T. Lin, W. Wang, *et al.*, "Planning-oriented autonomous driving," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 17853–17862, 2023.

[197] C. Cui, Y. Ma, X. Cao, W. Ye, Y. Zhou, K. Liang, J. Chen, J. Lu, Z. Yang, K.-D. Liao, *et al.*, "A survey on multimodal large language models for autonomous driving," *arXiv preprint arXiv:2311.12320*, 2023.

[198] X. Zhou, M. Liu, B. L. Zagar, E. Yurtsever, and A. C. Knoll, "Vision language models in autonomous driving and intelligent transportation systems," *arXiv preprint arXiv:2310.14414*, 2023.

[199] F.-Y. Wang, Q. Miao, L. Li, Q. Ni, X. Li, J. Li, L. Fan, Y. Tian, and Q.-L. Han, "When does sora show: The beginning of tao to imaginative intelligence and scenarios engineering," *IEEE/CAA Journal of Automatica Sinica*, vol. 11, no. 4, pp. 809–815, 2024.

[200] X. Li, Y. Tian, P. Ye, H. Duan, and F.-Y. Wang, "A novel scenarios engineering methodology for foundation models in metaverse," *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, vol. 53, no. 4, pp. 2148–2159, 2022.

[201] Y. Tian, X. Li, H. Zhang, C. Zhao, B. Li, X. Wang, and F.-Y. Wang, "Vistagpt: Generative parallel transformers for vehicles with intelligent systems for transport automation," *IEEE Transactions on Intelligent Vehicles*, 2023.

[202] S. W. Kim, J. Philion, A. Torralba, and S. Fidler, "Drivegan: Towards a controllable high-quality neural simulation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5820–5829, 2021.

[203] 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.

[204] 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.

[205] F. Jia, W. Mao, Y. Liu, Y. Zhao, Y. Wen, C. Zhang, X. Zhang, and T. Wang, "Adriver-i: A general world model for autonomous driving," *arXiv preprint arXiv:2311.13549*, 2023.

[206] T. Brooks, B. Peebles, C. Homes, W. DePue, Y. Guo, L. Jing, D. Schnurr, J. Taylor, T. Luhman, E. Luhman, C. Ng, R. Wang, and A. Ramesh, "Video generation models as world simulators," 2024.

[207] M. Schwonberg, J. Niemeijer, J.-A. Termöhlen, J. P. Schäfer, N. M. Schmidt, H. Gottschalk, and T. Fingscheidt, "Survey on unsupervised domain adaptation for semantic segmentation for visual perception in automated driving," *IEEE Access*, 2023.

[208] Ö. Erkent and C. Laugier, "Semantic segmentation with unsupervised domain adaptation under varying weather conditions for autonomous

vehicles," *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 3580–3587, 2020.

[209] Y. Wei, Z. Wei, Y. Rao, J. Li, J. Zhou, and J. Lu, "Lidar distillation: Bridging the beam-induced domain gap for 3d object detection," in *European Conference on Computer Vision*, pp. 179–195, Springer, 2022.

[210] C. Hu, S. Hudson, M. Ethier, M. Al-Sharman, D. Rayside, and W. Melek, "Sim-to-real domain adaptation for lane detection and classification in autonomous driving," in *2022 IEEE Intelligent Vehicles Symposium (IV)*, pp. 457–463, IEEE, 2022.

[211] E. Hüllermeier and W. Waegeman, "Aleatoric and epistemic uncertainty in machine learning: An introduction to concepts and methods," *Machine Learning*, vol. 110, pp. 457–506, 2021.

[212] S. Shafaei, S. Kugele, M. H. Osman, and A. Knoll, "Uncertainty in machine learning: A safety perspective on autonomous driving," in *Computer Safety, Reliability, and Security: SAFECOMP 2018 Workshops, ASSURE, DECSoS, SASSUR, STRIVE, and WAISE, Västerås, Sweden, September 18, 2018, Proceedings 37*, pp. 458–464, Springer, 2018.

[213] M. S. Gal and D. L. Rubinfeld, "Data standardization," *NYUL Rev.*, vol. 94, p. 737, 2019.

[214] A. Ess, B. Leibe, and L. Van Gool, "Depth and appearance for mobile scene analysis," in *2007 IEEE 11th international conference on computer vision*, pp. 1–8, IEEE, 2007.

[215] C. Wojek, S. Walk, and B. Schiele, "Multi-cue onboard pedestrian detection," in *2009 IEEE conference on computer vision and pattern recognition*, pp. 794–801, IEEE, 2009.

[216] W. Choi, K. Shahid, and S. Savarese, "What are they doing?: Collective activity classification using spatio-temporal relationship among people," in *2009 IEEE 12th international conference on computer vision workshops, ICCV Workshops*, pp. 1282–1289, IEEE, 2009.

[217] D. M. Chen, G. Baatz, K. Köser, S. S. Tsai, R. Vedantham, T. Pylvänäinen, K. Roimela, X. Chen, J. Bach, M. Pollefeys, *et al.*, "City-scale landmark identification on mobile devices," in *CVPR 2011*, pp. 737–744, IEEE, 2011.

[218] X. Li, F. Flohr, Y. Yang, H. Xiong, M. Braun, S. Pan, K. Li, and D. M. Gavrila, "A new benchmark for vision-based cyclist detection," in *2016 IEEE Intelligent Vehicles Symposium (IV)*, pp. 1028–1033, IEEE, 2016.

[219] R. Timofte, K. Zimmermann, and L. Van Gool, "Multi-view traffic sign detection, recognition, and 3d localisation," *Machine vision and applications*, vol. 25, pp. 633–647, 2014.

[220] A. Teichman, J. Levinson, and S. Thrun, "Towards 3d object recognition via classification of arbitrary object tracks," in *2011 IEEE International Conference on Robotics and Automation*, pp. 4034–4041, IEEE, 2011.

[221] C. Caraffi, T. Vojíř, J. Trefnỳ, J. Šochman, and J. Matas, "A system for real-time detection and tracking of vehicles from a single car-mounted camera," in *2012 15th international IEEE conference on intelligent transportation systems*, pp. 975–982, IEEE, 2012.

[222] J.-L. Blanco-Claraco, F.-A. Moreno-Duenas, and J. González-Jiménez, "The málaga urban dataset: High-rate stereo and lidar in a realistic urban scenario," *The International Journal of Robotics Research*, vol. 33, no. 2, pp. 207–214, 2014.

[223] M. De Deuge, A. Quadros, C. Hung, and B. Douillard, "Unsupervised feature learning for classification of outdoor 3d scans," in *Australasian conference on robitics and automation*, vol. 2, University of New South Wales Kensington, Australia, 2013.

[224] D. Pfeiffer, S. Gehrig, and N. Schneider, "Exploiting the power of stereo confidences," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 297–304, 2013.

[225] K. Matzen and N. Snavely, "Nyc3dcars: A dataset of 3d vehicles in geographic context," in *Proceedings of the IEEE International Conference on Computer Vision*, pp. 761–768, 2013.

[226] P. Koschorrek, T. Piccini, P. Oberg, M. Felsberg, L. Nielsen, and R. Mester, "A multi-sensor traffic scene dataset with omnidirectional video," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 727–734, 2013.

[227] N. Schneider and D. M. Gavrila, "Pedestrian path prediction with recursive bayesian filters: A comparative study," in *german conference on pattern recognition*, pp. 174–183, Springer, 2013.

[228] A. Møgelmose, D. Liu, and M. M. Trivedi, "Traffic sign detection for us roads: Remaining challenges and a case for tracking," in *17th International IEEE Conference on Intelligent Transportation Systems (ITSC)*, pp. 1394–1399, IEEE, 2014.

[229] A. Serna, B. Marcotegui, F. Goulette, and J.-E. Deschaud, "Paris-rue-madame database: a 3d mobile laser scanner dataset for benchmarking urban detection, segmentation and classification methods," in *4th international conference on pattern recognition, applications and methods ICPRAM 2014*, 2014.

[230] R. Guerrero-Gómez-Olmedo, B. Torre-Jiménez, R. López-Sastre, S. Maldonado-Bascón, and D. Onoro-Rubio, "Extremely overlapping vehicle counting," in *Pattern Recognition and Image Analysis: 7th Iberian Conference, IbPRIA 2015, Santiago de Compostela, Spain, June 17-19, 2015, Proceedings 7*, pp. 423–431, Springer, 2015.

[231] N. Mayer, E. Ilg, P. Hausser, P. Fischer, D. Cremers, A. Dosovitskiy, and T. Brox, "A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4040–4048, 2016.

[232] L. Wen, D. Du, Z. Cai, Z. Lei, M.-C. Chang, H. Qi, J. Lim, M.-H. Yang, and S. Lyu, "Ua-detrac: A new benchmark and protocol for multi-object detection and tracking," *Computer Vision and Image Understanding*, vol. 193, p. 102907, 2020.

[233] N. Carlevaris-Bianco, A. K. Ushani, and R. M. Eustice, "University of michigan north campus long-term vision and lidar dataset," *The International Journal of Robotics Research*, vol. 35, no. 9, pp. 1023–1035, 2016.

[234] R. Guzmán, J.-B. Hayet, and R. Klette, "Towards ubiquitous autonomous driving: The ccsad dataset," in *Computer Analysis of Images and Patterns: 16th International Conference, CAIP 2015, Valletta, Malta, September 2-4, 2015 Proceedings, Part I 16*, pp. 582–593, Springer, 2015.

[235] S. Hwang, J. Park, N. Kim, Y. Choi, and I. So Kweon, "Multispectral pedestrian detection: Benchmark dataset and baseline," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1037–1045, 2015.

[236] "Sap: https://cs.stanford.edu/ anenberg/ua_data/,"

[237] P. Pinggera, S. Ramos, S. Gehrig, U. Franke, C. Rother, and R. Mester, "Lost and found: detecting small road hazards for self-driving vehicles. in 2016 ieee," in *RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 1099–1106.

[238] E. Romera, L. M. Bergasa, and R. Arroyo, "Need data for driver behaviour analysis? presenting the public uah-driveset," in *2016 IEEE 19th international conference on intelligent transportation systems (ITSC)*, pp. 387–392, IEEE, 2016.

[239] D. Temel, G. Kwon, M. Prabhushankar, and G. AlRegib, "Cure-tsr: Challenging unreal and real environments for traffic sign recognition," *arXiv preprint arXiv:1712.02463*, 2017.

[240] "Tusimple: https://github.com/tusimple/tusimple-benchmark,"

[241] X. Liu, Z. Deng, H. Lu, and L. Cao, "Benchmark for road marking detection: Dataset specification and performance baseline," in *2017 IEEE 20th International Conference on Intelligent Transportation Systems (ITSC)*, pp. 1–6, IEEE, 2017.

[242] I. Klein, "Nexet—the largest and most diverse road dataset in the world," *Medium*, 2017.

[243] "Diml: https://dimlrgbd.github.io/,"

[244] K. Behrendt, L. Novak, and R. Botros, "A deep learning approach to traffic lights: Detection, tracking, and classification," in *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 1370–1377, IEEE, 2017.

[245] J. Jeong, Y. Cho, Y.-S. Shin, H. Roh, and A. Kim, "Complex urban dataset with multi-level sensors from highly diverse urban environments," *The International Journal of Robotics Research*, vol. 38, no. 6, pp. 642–657, 2019.

[246] Y. Hu, J. Binas, D. Neil, S.-C. Liu, and T. Delbruck, "Ddd20 end-to-end event camera driving dataset: Fusing frames and events with deep learning for improved steering prediction," in *2020 IEEE 23rd International Conference on Intelligent Transportation Systems (ITSC)*, pp. 1–6, IEEE, 2020.

[247] W. Kim, M. S. Ramanagopal, C. Barto, M.-Y. Yu, K. Rosaen, N. Goumas, R. Vasudevan, and M. Johnson-Roberson, "Pedx: Benchmark dataset for metric 3-d pose estimation of pedestrians in complex urban intersections," *IEEE Robotics and Automation Letters*, vol. 4, no. 2, pp. 1940–1947, 2019.

[248] Y. Chen, J. Wang, J. Li, C. Lu, Z. Luo, H. Xue, and C. Wang, "Lidar-video driving dataset: Learning driving policies effectively," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5870–5878, 2018.

[249] M. Wrenninge and J. Unger, "Synscapes: A photorealistic synthetic dataset for street scene parsing," *arXiv preprint arXiv:1810.08705*, 2018.

[250] A. Z. Zhu, D. Thakur, T. Özaslan, B. Pfrommer, V. Kumar, and K. Daniilidis, "The multivehicle stereo event camera dataset: An event

camera dataset for 3d perception," *IEEE Robotics and Automation Letters*, vol. 3, no. 3, pp. 2032–2039, 2018.

[251] T. Sattler, W. Maddern, C. Toft, A. Torii, L. Hammarstrand, E. Stenborg, D. Safari, M. Okutomi, M. Pollefeys, J. Sivic, *et al.*, "Benchmarking 6dof outdoor visual localization in changing conditions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 8601–8610, 2018.

[252] A. P. Shah, J.-B. Lamare, T. Nguyen-Anh, and A. Hauptmann, "Cadp: A novel dataset for cctv traffic camera based accident analysis," in *2018 15th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS)*, pp. 1–9, IEEE, 2018.

[253] T. Fleck, K. Daaboul, M. Weber, P. Schörner, M. Wehmer, J. Doll, S. Orf, N. Sußmann, C. Hubschneider, M. R. Zofka, *et al.*, "Towards large scale urban traffic reference data: Smart infrastructure in the test area autonomous driving baden-württemberg," in *Intelligent Autonomous Systems 15: Proceedings of the 15th International Conference IAS-15*, pp. 964–982, Springer, 2019.

[254] H. Schafer, E. Santana, A. Haden, and R. Biasini, "A commute in data: The comma2k19 dataset," *arXiv preprint arXiv:1812.05752*, 2018.

[255] D. Dai and L. Van Gool, "Dark model adaptation: Semantic image segmentation from daytime to nighttime," in *2018 21st International Conference on Intelligent Transportation Systems (ITSC)*, pp. 3819–3824, IEEE, 2018.

[256] X. Roynard, J.-E. Deschaud, and F. Goulette, "Paris-lille-3d: A large and high-quality ground-truth urban point cloud dataset for automatic segmentation and classification," *The International Journal of Robotics Research*, vol. 37, no. 6, pp. 545–557, 2018.

[257] L. Neumann, M. Karg, S. Zhang, C. Scharfenberger, E. Piegert, S. Mistr, O. Prokofyeva, R. Thiel, A. Vedaldi, A. Zisserman, *et al.*, "Nightowls: A pedestrians at night dataset," in *Computer Vision–ACCV 2018: 14th Asian Conference on Computer Vision, Perth, Australia, December 2–6, 2018, Revised Selected Papers, Part I 14*, pp. 691–705, Springer, 2019.

[258] M. Braun, S. Krebs, F. Flohr, and D. Gavrila, "The eurocity persons dataset: A novel benchmark for object detection. arxiv 2018," *arXiv preprint arXiv:1805.07193*.

[259] G. Choe, S.-H. Kim, S. Im, J.-Y. Lee, S. G. Narasimhan, and I. S. Kweon, "Ranus: Rgb and nir urban scene dataset for deep scene parsing," *IEEE Robotics and Automation Letters*, vol. 3, no. 3, pp. 1808–1815, 2018.

[260] D. Griffiths and J. Boehm, "Synthcity: A large scale synthetic point cloud," *arXiv preprint arXiv:1907.04758*, 2019.

[261] Z. Che, G. Li, T. Li, B. Jiang, X. Shi, X. Zhang, Y. Lu, G. Wu, Y. Liu, and J. Ye, "D$^2$-city: a large-scale dashcam video dataset of diverse traffic scenarios," *arXiv preprint arXiv:1904.01975*, 2019.

[262] M. Aly, "Real time detection of lane markers in urban streets," in *2008 IEEE intelligent vehicles symposium*, pp. 7–12, IEEE, 2008.

[263] Y. Dong, Y. Zhong, W. Yu, M. Zhu, P. Lu, Y. Fang, J. Hong, and H. Peng, "Mcity data collection for automated vehicles study," *arXiv preprint arXiv:1912.06258*, 2019.

[264] W. Cheng, H. Luo, W. Yang, L. Yu, S. Chen, and W. Li, "Det: A high-resolution dvs dataset for lane extraction," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pp. 0–0, 2019.

[265] A. Patil, S. Malla, H. Gang, and Y.-T. Chen, "The h3d dataset for full-surround 3d multi-object detection and tracking in crowded urban scenes," in *2019 International Conference on Robotics and Automation (ICRA)*, pp. 9552–9557, IEEE, 2019.

[266] K. Behrendt and R. Soussan, "Unsupervised labeled lane markers using maps," in *Proceedings of the IEEE/CVF international conference on computer vision workshops*, pp. 0–0, 2019.

[267] L. Ding, J. Terwilliger, R. Sherony, B. Reimer, and L. Fridman, "Value of temporal dynamics information in driving scene segmentation," *IEEE Transactions on Intelligent Vehicles*, vol. 7, no. 1, pp. 113–122, 2021.

[268] M. Meyer and G. Kuschk, "Automotive radar dataset for deep learning based 3d object detection," in *2019 16th european radar conference (EuRAD)*, pp. 129–132, IEEE, 2019.

[269] S. Nag, S. Adak, and S. Das, "What's there in the dark," in *2019 IEEE International Conference on Image Processing (ICIP)*, pp. 2996–3000, IEEE, 2019.

[270] K. Behrendt, "Boxy vehicle detection in large images," in *Proceedings of the IEEE/CVF international conference on computer vision workshops*, pp. 0–0, 2019.

[271] M. Wigness, S. Eum, J. G. Rogers, D. Han, and H. Kwon, "A rugd dataset for autonomous navigation and visual perception in unstructured outdoor environments," in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 5000–5007, IEEE, 2019.

[272] X. Song, P. Wang, D. Zhou, R. Zhu, C. Guan, Y. Dai, H. Su, H. Li, and R. Yang, "Apollocar3d: A large 3d car instance understanding benchmark for autonomous driving," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5452–5462, 2019.

[273] Y. Yao, M. Xu, C. Choi, D. J. Crandall, E. M. Atkins, and B. Dariush, "Egocentric vision-based future vehicle localization for intelligent driving assistance systems," in *2019 International Conference on Robotics and Automation (ICRA)*, pp. 9711–9717, IEEE, 2019.

[274] 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/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10591–10599, 2019.

[275] F. Codevilla, E. Santana, A. M. López, and A. Gaidon, "Exploring the limitations of behavior cloning for autonomous driving," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 9329–9338, 2019.

[276] A. Ligocki, A. Jelinek, and L. Zalud, "Brno urban dataset-the new data for self-driving agents and mapping tasks," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 3284–3290, IEEE, 2020.

[277] Y. Lou, Y. Bai, J. Liu, S. Wang, and L. Duan, "Veri-wild: A large dataset and a new method for vehicle re-identification in the wild," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 3235–3243, 2019.

[278] 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.

[279] K. Cordes and H. Broszio, "Vehicle lane merge visual benchmark," in *2020 25th International Conference on Pattern Recognition (ICPR)*, pp. 715–722, IEEE, 2021.

[280] A. Singh, A. Kamireddypalli, V. Gandhi, and K. M. Krishna, "Lidar guided small obstacle segmentation," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 8513–8520, IEEE, 2020.

[281] Z. Wang, S. Ding, Y. Li, J. Fenn, S. Roychowdhury, A. Wallin, L. Martin, S. Ryvola, G. Sapiro, and Q. Qiu, "Cirrus: A long-range bi-pattern lidar dataset," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 5744–5750, IEEE, 2021.

[282] E. Mohamed, M. Ewaisha, M. Siam, H. Rashed, S. Yogamani, W. Hamdy, M. El-Dakdouky, and A. El-Sallab, "Monocular instance motion segmentation for autonomous driving: Kitti instancemotseg dataset and multi-task baseline," in *2021 IEEE Intelligent Vehicles Symposium (IV)*, pp. 114–121, IEEE, 2021.

[283] Q.-H. Pham, P. Sevestre, R. S. Pahwa, H. Zhan, C. H. Pang, Y. Chen, A. Mustafa, V. Chandrasekhar, and J. Lin, "A 3d dataset: Towards autonomous driving in challenging environments," in *2020 IEEE International conference on Robotics and Automation (ICRA)*, pp. 2267–2273, IEEE, 2020.

[284] W. Tan, N. Qin, L. Ma, Y. Li, J. Du, G. Cai, K. Yang, and J. Li, "Toronto-3d: A large-scale mobile lidar dataset for semantic segmentation of urban roadways," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pp. 202–203, 2020.

[285] L. Ding, M. Glazer, M. Wang, B. Mehler, B. Reimer, and L. Fridman, "Mit-avt clustered driving scene dataset: Evaluating perception systems in real-world naturalistic driving scenarios," in *2020 IEEE Intelligent Vehicles Symposium (IV)*, pp. 232–237, IEEE, 2020.

[286] Y. Pan, B. Gao, J. Mei, S. Geng, C. Li, and H. Zhao, "Semanticposs: A point cloud dataset with large quantity of dynamic instances," in *2020 IEEE Intelligent Vehicles Symposium (IV)*, pp. 687–693, IEEE, 2020.

[287] E. Alberti, A. Tavera, C. Masone, and B. Caputo, "Idda: A large-scale multi-domain dataset for autonomous driving," *IEEE Robotics and Automation Letters*, vol. 5, no. 4, pp. 5526–5533, 2020.

[288] A. Ouaknine, A. Newson, J. Rebut, F. Tupin, and P. Pérez, "Carrada dataset: Camera and automotive radar with range-angle-doppler annotations," in *2020 25th International Conference on Pattern Recognition (ICPR)*, pp. 5068–5075, IEEE, 2021.

[289] S. Malla, B. Dariush, and C. Choi, "Titan: Future forecast using action priors," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11186–11196, 2020.

[290] Z. Xie, S. Wang, K. Xu, Z. Zhang, X. Tan, Y. Xie, and L. Ma, "Boosting night-time scene parsing with learnable frequency," *IEEE Transactions on Image Processing*, 2023.

[291] V. Guizilini, R. Ambrus, S. Pillai, A. Raventos, and A. Gaidon, "3d packing for self-supervised monocular depth estimation," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 2485–2494, 2020.

[292] G. Kim, Y. S. Park, Y. Cho, J. Jeong, and A. Kim, "Mulran: Multimodal range dataset for urban place recognition," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 6246–6253, IEEE, 2020.

[293] Y. Zhou, G. Wan, S. Hou, L. Yu, G. Wang, X. Rui, and S. Song, "Da4ad: End-to-end deep attention-based visual localization for autonomous driving," in *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXVIII 16*, pp. 271–289, Springer, 2020.

[294] Z. Yan, L. Sun, T. Krajník, and Y. Ruichek, "Eu long-term dataset with multiple sensors for autonomous driving," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 10697–10704, IEEE, 2020.

[295] M. Ramezani, Y. Wang, M. Camurri, D. Wisth, M. Mattamala, and M. Fallon, "The newer college dataset: Handheld lidar, inertial and vision with ground truth," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 4353–4360, IEEE, 2020.

[296] W. Bao, Q. Yu, and Y. Kong, "Uncertainty-based traffic accident anticipation with spatio-temporal relational learning," in *Proceedings of the 28th ACM International Conference on Multimedia*, pp. 2682–2690, 2020.

[297] A. Carballo, J. Lambert, A. Monrroy, D. Wong, P. Narksri, Y. Kitsukawa, E. Takeuchi, S. Kato, and K. Takeda, "Libre: The multiple 3d lidar dataset," in *2020 IEEE Intelligent Vehicles Symposium (IV)*, pp. 1094–1101, IEEE, 2020.

[298] T. Gruber, F. Julca-Aguilar, M. Bijelic, and F. Heide, "Gated2depth: Real-time dense lidar from gated images," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 1506–1516, 2019.

[299] J. Wiederer, A. Bouazizi, U. Kressel, and V. Belagiannis, "Traffic control gesture recognition for autonomous vehicles," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 10676–10683, IEEE, 2020.

[300] P. Wenzel, R. Wang, N. Yang, Q. Cheng, Q. Khan, L. von Stumberg, N. Zeller, and D. Cremers, "4seasons: A cross-season dataset for multi-weather slam in autonomous driving," in *Pattern Recognition: 42nd DAGM German Conference, DAGM GCPR 2020, Tübingen, Germany, September 28–October 1, 2020, Proceedings 42*, pp. 404–417, Springer, 2021.

[301] S. Saralajew, L. Ohnemus, L. Ewecker, E. Asan, S. Isele, and S. Roos, "A dataset for provident vehicle detection at night," in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 9750–9757, IEEE, 2021.

[302] H. Sakashita, C. Flothow, N. Takemura, and Y. Sugano, "Driv100: In-the-wild multi-domain dataset and evaluation for real-world domain adaptation of semantic segmentation," *arXiv preprint arXiv:2102.00150*, 2021.

[303] L. Wang, L. Lei, H. Song, and W. Wang, "The neolix open dataset for autonomous driving," *arXiv preprint arXiv:2011.13528*, 2020.

[304] P. Spannaus, P. Zechel, and K. Lenz, "Automatum data: Drone-based highway dataset for the development and validation of automated driving software for research and commercial applications," in *2021 IEEE Intelligent Vehicles Symposium (IV)*, pp. 1372–1377, IEEE, 2021.

[305] L. Li, K. N. Ismail, H. P. Shum, and T. P. Breckon, "Durlar: A high-fidelity 128-channel lidar dataset with panoramic ambient and reflectivity imagery for multi-modal autonomous driving applications," in *2021 International Conference on 3D Vision (3DV)*, pp. 1227–1237, IEEE, 2021.

[306] B. Helou, A. Dusi, A. Collin, N. Mehdipour, Z. Chen, C. Lizarazo, C. Belta, T. Wongpiromsarn, R. D. Tebbens, and O. Beijbom, "The reasonable crowd: Towards evidence-based and interpretable models of driving behavior," in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 6708–6715, IEEE, 2021.

[307] M. Fabbri, G. Brasó, G. Maugeri, O. Cetintas, R. Gasparini, A. Ošep, S. Calderara, L. Leal-Taixé, and R. Cucchiara, "Motsynth: How can synthetic data help pedestrian detection and tracking?," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 10849–10859, 2021.

[308] F. R. Valverde, J. V. Hurtado, and A. Valada, "There is more than meets the eye: Self-supervised multi-object detection and tracking with sound

[309] by distilling multimodal knowledge," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11612–11621, 2021.

[309] N. A. M. Mai, P. Duthon, L. Khoudour, A. Crouzil, and S. A. Velastin, "3d object detection with sls-fusion network in foggy weather conditions," *Sensors*, vol. 21, no. 20, p. 6711, 2021.

[310] Y. Yuan and M. Sester, "Comap: A synthetic dataset for collective multi-agent perception of autonomous driving," *The International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, vol. 43, pp. 255–263, 2021.

[311] J. Oh, G. Lee, J. Park, W. Oh, J. Heo, H. Chung, D. H. Kim, B. Park, C.-G. Lee, S. Choi, *et al.*, "Towards defensive autonomous driving: Collecting and probing driving demonstrations of mixed qualities," in *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 12528–12533, IEEE, 2022.

[312] M. Howe, I. Reid, and J. Mackenzie, "Weakly supervised training of monocular 3d object detectors using wide baseline multi-view traffic camera data," *arXiv preprint arXiv:2110.10966*, 2021.

[313] O. Jayasinghe, S. Hemachandra, D. Anhettigama, S. Kariyawasam, R. Rodrigo, and P. Jayasekara, "Ceymo: see more on roads-a novel benchmark dataset for road marking detection," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 3104–3113, 2022.

[314] J. Jin, A. Fatemi, W. M. P. Lira, F. Yu, B. Leng, R. Ma, A. Mahdavi-Amiri, and H. Zhang, "Raidar: A rich annotated image dataset of rainy street scenes," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 2951–2961, 2021.

[315] H. Blum, P.-E. Sarlin, J. Nieto, R. Siegwart, and C. Cadena, "The fishyscapes benchmark: Measuring blind spots in semantic segmentation," *International Journal of Computer Vision*, vol. 129, pp. 3119–3135, 2021.

[316] O. Schumann, M. Hahn, N. Scheiner, F. Weishaupt, J. F. Tilly, J. Dickmann, and C. Wöhler, "Radarscenes: A real-world radar point cloud data set for automotive applications," in *2021 IEEE 24th International Conference on Information Fusion (FUSION)*, pp. 1–8, IEEE, 2021.

[317] G. Singh, S. Akrigg, M. Di Maio, V. Fontana, R. J. Alitappeh, S. Khan, S. Saha, K. Jeddisaravi, F. Yousefi, J. Culley, *et al.*, "Road: The road event awareness dataset for autonomous driving," *IEEE transactions on pattern analysis and machine intelligence*, vol. 45, no. 1, pp. 1036–1054, 2022.

[318] J. Han, X. Liang, H. Xu, K. Chen, L. Hong, J. Mao, C. Ye, W. Zhang, Z. Li, X. Liang, *et al.*, "Soda10m: a large-scale 2d self/semi-supervised object detection dataset for autonomous driving," *arXiv preprint arXiv:2106.11118*, 2021.

[319] J.-L. Déziel, P. Merriaux, F. Tremblay, D. Lessard, D. Plourde, J. Stanguennec, P. Goulet, and P. Olivier, "Pixset: An opportunity for 3d computer vision to go beyond point clouds with a full-waveform lidar dataset," in *2021 IEEE International Intelligent Transportation Systems Conference (ITSC)*, pp. 2987–2993, IEEE, 2021.

[320] R. Chan, K. Lis, S. Uhlemeyer, H. Blum, S. Honari, R. Siegwart, P. Fua, M. Salzmann, and M. Rottmann, "Segmentmeifyoucan: A benchmark for anomaly segmentation," *arXiv preprint arXiv:2104.14812*, 2021.

[321] Y. Zhang, L. Zhu, W. Feng, H. Fu, M. Wang, Q. Li, C. Li, and S. Wang, "Vil-100: A new dataset and a baseline model for video instance lane detection," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 15681–15690, 2021.

[322] X. Zhang, Z. Li, Y. Gong, D. Jin, J. Li, L. Wang, Y. Zhu, and H. Liu, "Openmpd: An open multimodal perception dataset for autonomous driving," *IEEE Transactions on Vehicular Technology*, vol. 71, no. 3, pp. 2437–2447, 2022.

[323] A. Kurup and J. Bos, "Dsor: A scalable statistical filter for removing falling snow from lidar point clouds in severe winter weather," *arXiv preprint arXiv:2109.07078*, 2021.

[324] J. Zhang, X. Zou, L.-D. Kuang, J. Wang, R. S. Sherratt, and X. Yu, "Cctsdb 2021: a more comprehensive traffic sign detection benchmark," *Human-centric Computing and Information Sciences*, vol. 12, 2022.

[325] P. Jiang and S. Saripalli, "Lidarnet: A boundary-aware domain adaptation model for point cloud semantic segmentation," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 2457–2464, IEEE, 2021.

[326] J. Gebele, B. Stuhr, and J. Haselberger, "Carlane: A lane detection benchmark for unsupervised domain adaptation from simulation to multiple real-world domains," *arXiv preprint arXiv:2206.08083*, 2022.

[327] A. Lehner, S. Gasperini, A. Marcos-Ramiro, M. Schmidt, M.-A. N. Mahani, N. Navab, B. Busam, and F. Tombari, "3d-vfield: Adversarial augmentation of point clouds for domain generalization in 3d object

detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 17295–17304, 2022.

[328] A. Kloukiniotis, A. Papandreou, C. Anagnostopoulos, A. Lalos, P. Kapsalas, D.-V. Nguyen, and K. Moustakas, "Carlascenes: A synthetic dataset for odometry in autonomous driving," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4520–4528, 2022.

[329] R. Xu, H. Xiang, X. Xia, X. Han, J. Li, and J. Ma, "Opv2v: An open benchmark dataset and fusion pipeline for perception with vehicle-to-vehicle communication," in *2022 International Conference on Robotics and Automation (ICRA)*, pp. 2583–2589, IEEE, 2022.

[330] K. Maag, R. Chan, S. Uhlemeyer, K. Kowol, and H. Gottschalk, "Two video data sets for tracking and retrieval of out of distribution objects," in *Proceedings of the Asian Conference on Computer Vision*, pp. 3776–3794, 2022.

[331] K. Cordes, C. Reinders, P. Hindricks, J. Lammers, B. Rosenhahn, and H. Broszio, "Roadsaw: A large-scale dataset for camera-based road surface and wetness estimation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4440–4449, 2022.

[332] H. Quispe, J. Sumire, P. Condori, E. Alvarez, and H. Vera, "I see you: A vehicle-pedestrian interaction dataset from traffic surveillance cameras," *arXiv preprint arXiv:2211.09342*, 2022.

[333] X. Wang, Z. Zhu, Y. Zhang, G. Huang, Y. Ye, W. Xu, Z. Chen, and X. Wang, "Are we ready for vision-centric driving streaming perception? the asap benchmark," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9600–9610, 2023.

[334] J. Breitenstein and T. Fingscheidt, "Amodal cityscapes: a new dataset, its generation, and an amodal semantic segmentation challenge baseline," in *2022 IEEE Intelligent Vehicles Symposium (IV)*, pp. 1018–1025, IEEE, 2022.

[335] A. R. Sekkat, Y. Dupuis, V. R. Kumar, H. Rashed, S. Yogamani, P. Vasseur, and P. Honeine, "Synwoodscape: Synthetic surround-view fisheye camera dataset for autonomous driving," *IEEE Robotics and Automation Letters*, vol. 7, no. 3, pp. 8502–8509, 2022.

[336] L. Zheng, Z. Ma, X. Zhu, B. Tan, S. Li, K. Long, W. Sun, S. Chen, L. Zhang, M. Wan, *et al.*, "Tj4dradset: A 4d radar dataset for autonomous driving," in *2022 IEEE 25th International Conference on Intelligent Transportation Systems (ITSC)*, pp. 493–498, IEEE, 2022.

[337] K. Li, K. Chen, H. Wang, L. Hong, C. Ye, J. Han, Y. Chen, W. Zhang, C. Xu, D.-Y. Yeung, *et al.*, "Coda: A real-world road corner case dataset for object detection in autonomous driving," in *European Conference on Computer Vision*, pp. 406–423, Springer, 2022.

[338] M. Hahner, C. Sakaridis, M. Bijelic, F. Heide, F. Yu, D. Dai, and L. Van Gool, "Lidar snowfall simulation for robust 3d object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16364–16374, 2022.

[339] G. Franchi, X. Yu, A. Bursuc, A. Tena, R. Kazmierczak, S. Dubuisson, E. Aldea, and D. Filliat, "Muad: Multiple uncertainties for autonomous driving, a benchmark for multiple uncertainty types and tasks," *arXiv preprint arXiv:2203.01437*, 2022.

[340] J. Cui, H. Qiu, D. Chen, P. Stone, and Y. Zhu, "Coopernaut: End-to-end driving with cooperative perception for networked vehicles," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 17252–17262, 2022.

[341] Z. Bai, G. Wu, M. J. Barth, Y. Liu, E. A. Sisbot, and K. Oguchi, "Pillargrid: Deep learning-based cooperative perception for 3d object detection from onboard-roadside lidar," in *2022 IEEE 25th International Conference on Intelligent Transportation Systems (ITSC)*, pp. 1743–1749, IEEE, 2022.

[342] D.-H. Paek, S.-H. Kong, and K. T. Wijaya, "K-lane: Lidar lane dataset and benchmark for urban roads and highways," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4450–4459, 2022.

[343] C. A. Diaz-Ruiz, Y. Xia, Y. You, J. Nino, J. Chen, J. Monica, X. Chen, K. Luo, Y. Wang, M. Emond, *et al.*, "Ithaca365: Dataset and driving perception under repeated and challenging weather conditions," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 21383–21392, 2022.

[344] N. Gray, M. Moraes, J. Bian, A. Wang, A. Tian, K. Wilson, Y. Huang, H. Xiong, and Z. Guo, "Glare: A dataset for traffic sign detection in sun glare," *IEEE Transactions on Intelligent Transportation Systems*, 2023.

[345] J. Hou, Q. Chen, Y. Cheng, G. Chen, X. Xue, T. Zeng, and J. Pu, "Sups: A simulated underground parking scenario dataset for autonomous

driving," in *2022 IEEE 25th International Conference on Intelligent Transportation Systems (ITSC)*, pp. 2265–2271, IEEE, 2022.

[346] K. Burnett, D. J. Yoon, Y. Wu, A. Z. Li, H. Zhang, S. Lu, J. Qian, W.-K. Tseng, A. Lambert, K. Y. Leung, *et al.*, "Boreas: A multi-season autonomous driving dataset," *The International Journal of Robotics Research*, vol. 42, no. 1-2, pp. 33–42, 2023.

[347] L. Kong, Y. Liu, X. Li, R. Chen, W. Zhang, J. Ren, L. Pan, K. Chen, and Z. Liu, "Robo3d: Towards robust and reliable 3d perception against corruptions," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 19994–20006, 2023.

[348] D.-H. Paek, S.-H. Kong, and K. T. Wijaya, "K-radar: 4d radar object detection for autonomous driving in various weather conditions," *Advances in Neural Information Processing Systems*, vol. 35, pp. 3819–3829, 2022.

[349] T. Matuszka, I. Barton, Á. Butykai, P. Hajas, D. Kiss, D. Kovács, S. Kunsági-Máté, P. Lengyel, G. Németh, L. Pető, *et al.*, "aimotive dataset: A multimodal dataset for robust autonomous driving with long-range perception," *arXiv preprint arXiv:2211.09445*, 2022.

[350] M. Büchner, J. Zürn, I.-G. Todoran, A. Valada, and W. Burgard, "Learning and aggregating lane graphs for urban automated driving," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 13415–13424, 2023.

[351] A. Marathe, D. Ramanan, R. Walambe, and K. Kotecha, "Wedge: A multi-weather autonomous driving dataset built from generative vision-language models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3317–3326, 2023.

[352] H. Wang, T. Li, Y. Li, L. Chen, C. Sima, Z. Liu, B. Wang, P. Jia, Y. Wang, S. Jiang, *et al.*, "Openlane-v2: A topology reasoning benchmark for unified 3d hd mapping," in *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023.

[353] Y. Li, S. Li, X. Liu, M. Gong, K. Li, N. Chen, Z. Wang, Z. Li, T. Jiang, F. Yu, *et al.*, "Sscbench: A large-scale 3d semantic scene completion benchmark for autonomous driving," *arXiv preprint arXiv:2306.09001*, 2023.

[354] K. Cordes and H. Broszio, "Camera-based road snow coverage estimation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4011–4019, 2023.

[355] H. Xiang, Z. Zheng, X. Xia, R. Xu, L. Gao, Z. Zhou, X. Han, X. Ji, M. Li, Z. Meng, *et al.*, "V2x-real: a largs-scale dataset for vehicle-to-everything cooperative perception," *arXiv preprint arXiv:2403.16034*, 2024.

[356] R. Hao, S. Fan, Y. Dai, Z. Zhang, C. Li, Y. Wang, H. Yu, W. Yang, J. Yuan, and Z. Nie, "Rcooper: A real-world large-scale dataset for roadside cooperative perception," *arXiv preprint arXiv:2403.10145*, 2024.