

A Survey on Autonomous Driving Datasets: Data Statistic, Annotation, and Outlook

Mingyu Liu[✉], Ekim Yurtsever[✉], *Member, IEEE*, Xingcheng Zhou[✉], Jonathan Fossaert[✉], Yuning Cui[✉],
Bare Luka Zagar[✉], Alois C. Knoll[✉], *Fellow, IEEE*

Abstract—Autonomous driving has rapidly developed and shown promising performance with recent advances in hardware and deep learning methods. High-quality datasets are fundamental for developing reliable autonomous driving algorithms. Previous dataset surveys tried to review the datasets but either focused on a limited number or lacked detailed investigation of the characters of datasets. To this end, we present an exhaustive study of over 200 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 each dataset, which can also be a guide for establishing new datasets. We further analyze the annotation process and quality of datasets. Additionally, we conduct an in-depth analysis of the data distribution of several vital datasets. Finally, we discuss the development trend of the future autonomous driving datasets.

Index Terms—Dataset, influence, annotation, autonomous driving.

I. INTRODUCTION

AUTONOMOUS 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 robotaxi [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 the real world [13]–[28], to name just a few. Secondly, the datasets vary in composition, including but not limited to multiple sensory data (like camera images and LiDAR point clouds) and different annotation types for various tasks in autonomous driving. Fig. 1 depicts the statistic of the 3D

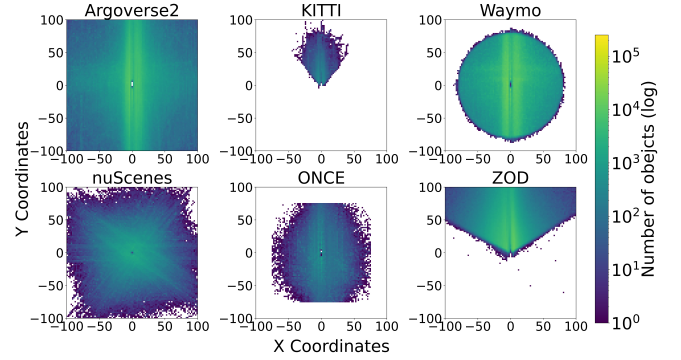


Fig. 1. BEV object distribution of datasets. Each heatmap corresponds to a dataset plotted with X and Y coordinates. Y is the driving direction of the ego-vehicle. The unique annotation characters of each dataset are reflected in the distribution range, density, and number of bounding boxes.

object bounding box distribution of six real-world datasets (Argoverse 2 [28], KITTI [13], nuScenes [22], ONCE [29], Waymo [23], and ZOD [30]) under a bird's-eye-view perspective, showing each dataset's distinguished annotation characters. According to the equipment position of sensors, the variety of datasets also reflects in the sensing domains, containing onboard, V2X, drone, and others. Furthermore, the geometrical diversity and weather conditions improve the generalization of the autonomous driving datasets.

A. Research Gap & Motivation

We demonstrate the yearly published number of perception datasets in Fig. 2 to reflect the trend of autonomous driving datasets from one perspective. Since numerous and continually increasing numbers of publicly published datasets exist, an exhaustive survey on autonomous driving datasets is valuable for advancing academic and industrial research. In prior work, Yin et al. [31] summarized 27 publicly available datasets containing data collected on public roads. As a sequential work of [31], [32] extended the dataset number. [33] and [34] proposed systematical introductions to the existing datasets from an application perspective. Beyond describing existing datasets, [35] discussed the domain adaptation between synthetic and real data and automatic labeling methods. [36] 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 [37] collected a large number of datasets while lacking detailed analysis for the attributes of

M. Liu, X. Zhou, Y. Cui, B.L. Zagar, and A.C. 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@tum.de, xingcheng.zhou@tum.de, yuning.cui@in.tum.de, bare.luka.zagar@tum.de, 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)

Survey	Year	General		S. domain	S. moda.	Data		Data analysis	Annotation	
		#Datasets	Tasks			Geo.	Env.		Quality	Process
When to use what dataset [31]	2017	27	Perc		✓	✓	✓			
Self-driving Algorithm [32]	2019	37	Perc		✓	✓	✓			
Is it safe to drive [33]	2019	54	Perc, Pred, E2E		✓	✓	✓			
CV for AVs [34]	2020	33	Perc		✓	✓	✓			✓
A Survey on AD Datasets [35]	2021	30	Perc		✓	✓	✓			✓
3D Semantic Segmentation [40]	2021	29	Perc		✓	✓	✓	✓	✓	✓
AD-Dataset [37]	2022	204	Perc, Pred, Pl, C	✓	✓	✓	✓			
Anomaly Detection [38]	2023	16	Perc		✓	✓	✓	✓		
Synthetic Datasets for AD [39]	2023	17	Perc, Pred		✓	✓	✓			
Decision-making [41]	2023	25	Pl, C		✓	✓	✓			
Open-sourced Data Ecosystem [36]	2023	70	Perc, Pred, Pl	✓	✓	✓	✓			✓
Ours		235	Perc, Pred, Pl, C, E2E	✓	✓	✓	✓	✓	✓	✓

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. WE DESCRIBE THE TASK TYPES IN A COARSE GRANULARITY, INCLUDING PERCEPTION (PERC.), PREDICTION (PRED.), PLANNING (PL.), CONTROL (C.) AND END-TO-END (E2E).

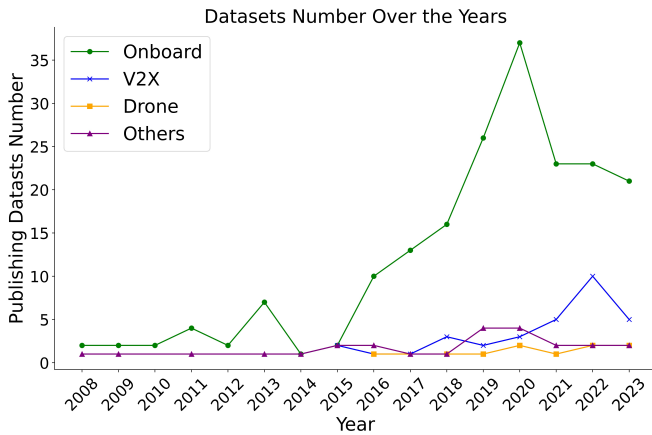


Fig. 2. Overview of the trend of publishing datasets. The diagram exhibits a rapid growth of onboard datasets between 2015 and 2020 and then slowly down. In contrast, there is an increment trend for the V2X datasets, showing the research trend on the cooperative perception systems.

these datasets. Compared to studies on all types of datasets, some researchers presented surveys on a particular type of autonomous driving dataset, such as anomaly detection [38], synthetic datasets [39], 3D semantic segmentation [40], and decision-making [41].

To this end, we aim to propose a comprehensive and systematic study on a large number of datasets in autonomous driving, covering all tasks from perception to control, considering real-world as well as synthetic data, and providing insight into the data modality and quality of several crucial datasets. We illustrate the comparison between other dataset surveys and ours in Tab I.

B. Main Contributions

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

- We present an exhaustive survey on autonomous driving datasets. We consider 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. To the best of our

knowledge, we provide an overview of the most extensive collection of autonomous driving datasets recorded to date.

- We systematically illustrate the sensors and sensing domains for collecting autonomous driving data. Furthermore, we describe the main tasks in autonomous driving, including task goals, required data modalities, and evaluation metrics.
- We summarize and divide the datasets according to their sensing domains and supporting tasks, which assists researchers in efficiently selecting and gathering information for the goal datasets. Thereby facilitating more targeted and effective research and development efforts.
- Additionally, we introduce an impact score metric to evaluate the influence of published perception datasets in the community. This metric can also be leveraged as a guidance for developing future datasets. We deeply analyze datasets with the highest grades based on the impact score, highlighting their advantages and utility.
- We investigate the annotation quality of datasets and the existing labeling procedures for various autonomous driving tasks.
- Detailed data statistics is conducted to demonstrate the data distribution of various datasets from different perspectives, exhibiting their inherent limitations and suitable use situations.
- We analyze the recent technology trend and demonstrate the development direction of the next generation of datasets. We also provide a prospect on the potential effect of the advancement of Large Language Models on future autonomous driving.

C. Scope & Limitations

We aimed to conduct an exhaustive survey on the existing autonomous driving datasets to provide assistance and guidance for developing future algorithms and datasets in the domain. We collected datasets focusing on the four basic autonomous driving tasks: perception, prediction, planning, and control. Because several versatile datasets support multiple tasks, we only explained them in the main scope they support to avoid repeat introduction. Additionally, we collected a large

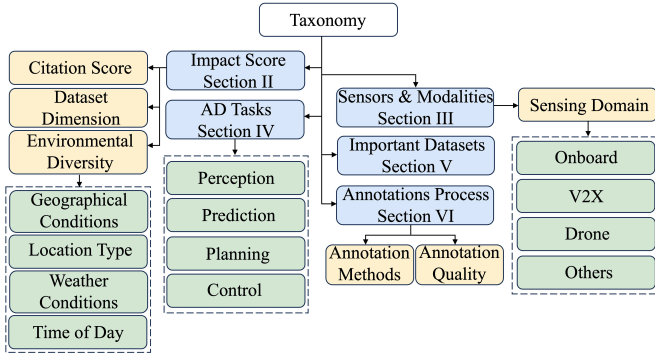


Fig. 3. This survey's primary taxonomy includes impact score, sensors and modalities, autonomous driving tasks, important datasets, and annotation process.

number of datasets and exhibited them with their primary characters in tables. Nevertheless, a detailed explanation of all collected datasets could not highlight the most popular ones, and it can hurt assisting researchers in finding valuable datasets through this survey. Hence, we only describe the most impactful datasets in detail.

D. Structure of the Article

The rest of the article is structured as follows: Section II introduces the approach we leveraged to obtain public datasets and the evaluation metrics for datasets. Section III demonstrates the primary sensors used in autonomous driving and their modalities. Section IV discusses autonomous driving tasks, related challenges, and required data. We further discuss several important datasets in Section V. The process of annotations and factors affecting annotation quality are exhibited in Section VI. Moreover, we statistic the data distribution of several datasets in section VII. In section VIII, we investigate the developing trend and future works of autonomous driving datasets. In the end, we conclude our work in Section IX. The taxonomy of this survey is shown in Fig. 3.

II. METHODOLOGY

This section consists of 1) How we collect and filter datasets (II-A), and 2) How to evaluate the impact of a dataset on the autonomous driving domain (II-B).

A. Datasets Collection

Following [42], we conducted a systematic review to exhaustively collect published autonomous driving datasets.

To ensure source diversity, we utilized well-known search engines such as Google¹, Google Scholar² and Baidu³ 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 datasets," "intelligent vehicle datasets,"

and terms related to object detection, classification, tracking, segmentation, prediction, planning, and control.

Furthermore, we explored IEEE Xplore⁴ 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 Github⁵ repositories and Paperwithcodes⁶. 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 publish time. Moreover, all citation counts were collected as of September 20, 2023, 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 the yearly citation increase speed of a dataset. The calculate 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 unbalance and highlight the differences of each dataset, we apply logarithmic transformation followed by Min-Max normalization to both c^t and c^a , described in Eq. 2.

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

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

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

¹<https://www.google.com/>

²<https://scholar.google.com/>

³<https://www.baidu.com/>

⁴<https://ieeexplore.ieee.org/Xplore/home.jsp>

⁵<https://github.com/>

⁶<https://paperswithcode.com/>

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_{normal} , 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, $l = \{1, 2, 3, 4, 5\}$, based on quintiles. To simplify the calculation, we normalize t_n and l before the following process.

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.

$$d_{\text{score}} = 0.5f_{\text{norm}} + 0.1t + 0.2t_{n\text{norm}} + 0.2l_{\text{norm}} \quad (4)$$

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 refers to when the data is collected, like morning or dusk. 3) the types of driving scenarios, e.g., urban or rural. 4) the geometric scope means the number of countries or cities where the data is recorded. It is worth noting that we treat the geometric scope for synthetic datasets as missing. 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 use the median value as the missing value. Otherwise, we set the missing value of this attribute 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} .

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

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. DATA SOURCES AND COOPERATIVE PERCEPTION IN AUTONOMOUS DRIVING

In this section, we introduce the sensors and their modalities mainly used in autonomous driving (III-A). Furthermore, in

III-B, we analyze the data acquisition and communication domain, such as onboard, drone, and vehicle-to-everything.

A. Sensors and Modalities of Data

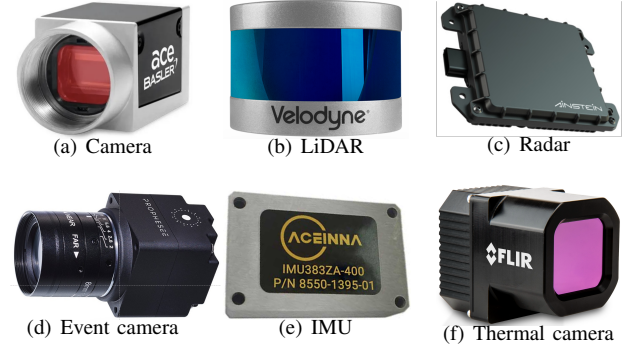


Fig. 4. Sensors on autonomous driving vehicle. The type of each sensor is (a) Camera: Basler ace acA1600-20uc, (b) LiDAR: Velodyne Puck LITE, (c) Radar: Ainstein Launches K-79, (d) Event-based camera: Evaluation Kit 4 HD, (e) IMU: IMU383_Aceinna-W and (f) Thermal camera: FLIR_2nd_Gen_ADK. All figures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

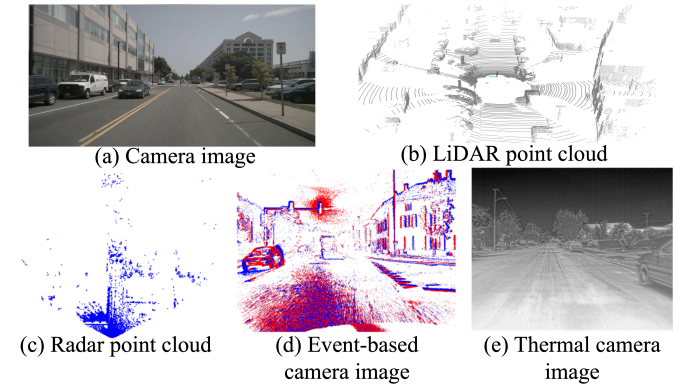


Fig. 5. Modalities of sensors. We demonstrate each type of sensor's modality to intuitively understand each sensor's characters.

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, LiDAR, and radars. Event-based and thermal cameras are also installed in vehicles or roadsides 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. All these cameras channel light through a lens onto an image sensor, e.g., CMOS, converting this light into electronic signals representing an image. 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 [43].

LiDAR Point Clouds. LiDAR uses laser beams to measure the distance between the sensor and an object, creating a 3D environment representation [44]. The 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. The weather conditions, e.g., fog, also limit the performance of LiDAR. In general, LiDAR is 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 [45]. 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 [46].

Infrared Images of Thermal Camera. Thermal cameras (see Fig. 5 (e)) detect heat signatures by capturing infrared radiation [47]. 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 [48]. In autonomous driving, it is used to track the movement and orientation of the vehicle. Although IMU does not include visual information of the surrounding environment, the perception systems can achieve more accurate and robust tracking of a vehicle's motion and orientation by fusing the data from an IMU with data from other sensors.

We analyze the sensor distribution from the collected datasets, shown in Fig. 6. More than half of the sensors are monocular cameras (53.85%) due to their low price and reliable performance. Additionally, 93 datasets 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, 5.41%, 3.42%, and 1.71% datasets include radar, thermal camera, and fisheye camera, respectively. Given the temporal efficiency of capturing dynamic scenes, three datasets generate data based on event-based cameras.

B. Sensing Domains and Cooperative Perception Systems

The sensory data and communication between the ego vehicle and other entities in the surrounding environment

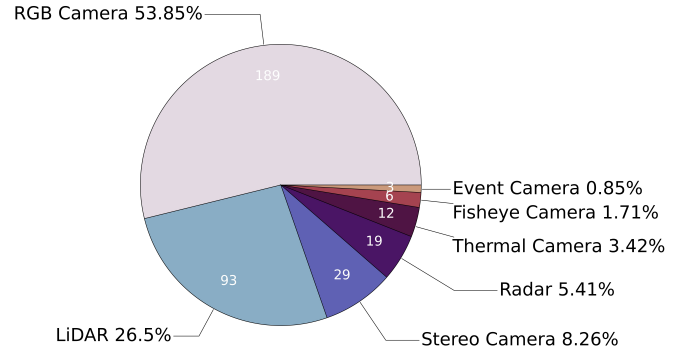


Fig. 6. Sensor number distribution. We statistic the distribution of different sensors. Overall, RGB cameras and LiDAR are the most used sensors in autonomous driving datasets.

play a pivotal role in ensuring self-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, sensors in the autonomous driving environment can be categorized into the following domains: onboard, Vehicle-to-Everything (V2X), drone, and others.

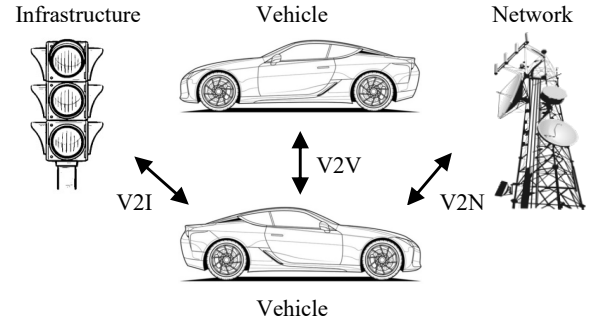


Fig. 7. Overview of cooperative perception systems in autonomous driving. A complete autonomous driving perception system consists of the ego vehicle and the cooperation between vehicles, infrastructure, and networks.

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

V2X. Vehicle-to-Everything (V2X) encompasses communications between a vehicle and any other components in the transport system [49], including V2V, V2I, and V2N (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 data, including their positions and velocity and sensory data, like camera

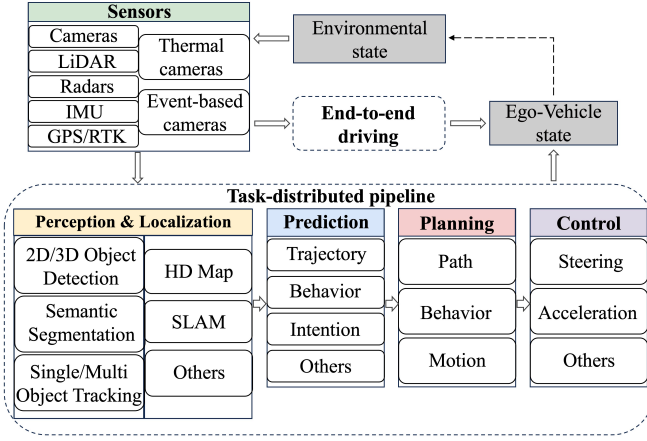


Fig. 8. The overview of autonomous driving pipeline. Autonomous driving systems can be categorized into two types: modular-based and end-to-end. Both rely on the data collected by various sensors installed on the vehicles or infrastructures and interact with the surrounding driving experiments.

images or LiDAR scans, which assists in creating a more comprehensive understanding of the 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 embedded in the road infrastructure, including cameras, LiDAR, radars, or event-based cameras, work synergistically to extend the perception range and improve the situational awareness of autonomous vehicles. In this survey, we consider both perceptions through infrastructure or V2I belonging to 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 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 illustrates examples of 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 and classify other entities within the driving environment. While 2D object detection identifies objects in image space, 3D object detection further incorporates precise depth information, often provided by LiDAR. 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 [54] is applied to evaluate the object detection performance. According to [1], the AP metric can be formulated as

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

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

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 [55], the main metrics used for segmentation are mean Pixel Accuracy (mPA):

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

and the mean Intersection over Union (mIoU):

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

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 [56] (see Eq. 9).

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

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 [57].

HD Map. HD mapping aims to construct detailed, highly accurate representations that include information about road structures, traffic signs, and landmarks. A dataset should

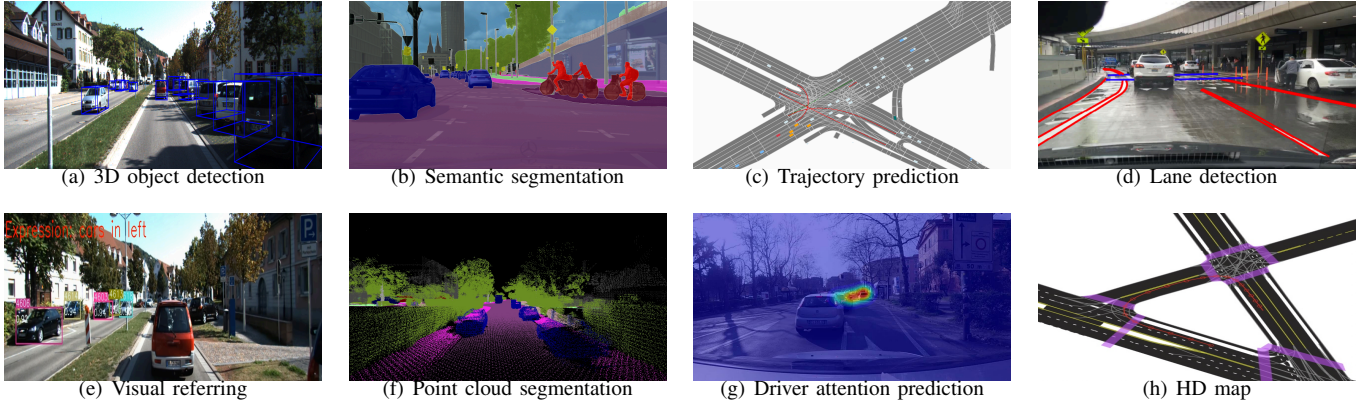


Fig. 9. Examples of various autonomous driving tasks. (a) is from KITTI [13], (b) is from Cityscapes [14], (c) is from V2X-Seq [50], (d) is from BDD100K [24], (e) is from Refer-KITTI [51], (f) is from KITTI-360 [52], (g) is from Dr(eye)ve [53], (h) is from Argoverse 2 [28]. All figures are collected from the open-source data of datasets or the websites hosting the datasets.

provide LiDAR data for precise spatial information and camera data for visual details to ensure established map accuracy. According to [28], HD map automation [58] and HD map change detection [59] have received more and more attention. Normally, 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. [60] 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 [61], [62], such as Root Mean Squared Error (RMSE):

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

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 [63].

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

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. Using time-series data from sensors like cameras and LiDAR, trajectory prediction pertains to

anticipating the future paths or movement patterns of other entities [61], such as pedestrians, cyclists, or other vehicles.

Behavior Prediction. Behavior prediction anticipates the potential actions of other road users [62], e.g., whether a vehicle will change lanes. 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 of the intention behind the actions of objects, involving a deeper semantic comprehension of the physical or mental activities of humans [64]. Because of the task's complexity, it requires data not only from perception sensors like cameras but also additional information, such as traffic signals and hand gestures, to infer the intentions of other agents.

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 consists of path, behavioral, and motion planning [53].

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 [53].

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 [65]. It translates high-level commands into actionable throttle, brake, and steering commands.

Dataset	Year	Size	Temp	2D Det	3D Det	2D Seg	3D Seg	Tracking	Lane Det	Categories number	Weather conditions	Time of day	Scenario type	Geometrical scope	Impact score
Onboard															
nuScenes [22]	2019	40K	✓	✓	✓	✓	✓	✓	✓	23	3	2	4	2	81.20
CyScapes [14]	2016	25K	×	✓	✓	✓	✓	✓	✓	30	1	1	3	1	79.00
BDD100K [24]	2020	12M	✓	✓	✓	✓	✓	✓	✓	40	5	2	4	1	77.82
Waymo [23]	2019	230K	✓	✓	✓	✓	✓	✓	✓	23	2	3	5	1	77.81
KITTI [13]	2012	41K	✓	✓	✓	✓	✓	✓	✓	8	1	1	2	1	77.30
SYNTHIA [2]	2016	13.4K	×	✓	✓	✓	✓	✓	✓	13	2	3	5	-	70.44
ApolloScapes [17]	2018	143,906	✓	✓	✓	✓	✓	✓	✓	28	4	3	3	1	68.97
Virtual KITTI [3]	2016	21,260	✓	✓	✓	✓	✓	✓	✓	8	3	2	5	-	67.61
VIPER [5]	2017	254,064	✓	✓	✓	✓	✓	✓	✓	11	4	3	4	-	67.07
SemanticKITTI [6]	2019	43,552	✓	✓	✓	✓	✓	✓	✓	28	1	1	4	1	66.66
GTA5 [4]	2016	24,966	×	✓	✓	✓	✓	✓	✓	19	2	2	2	-	66.10
Argoverse 2 [28]	2023	6M	✓	✓	✓	✓	✓	✓	✓	30	2	1	1	6	61.11
Lane Det [68]	2017	133,235	×	✓	✓	✓	✓	✓	✓	1	1	3	3	1	59.62
CityPersons [69]	2017	25K	×	✓	✓	✓	✓	✓	✓	30	2	1	-	27	58.84
CamVid [70]	2009	701	×	✓	✓	✓	✓	✓	✓	32	2	1	2	1	58.69
IDD [71]	2019	10,004	×	✓	✓	✓	✓	✓	✓	34	3	3	5	1	58.12
Foggy Cityscapes [72]	2018	20,550	×	✓	✓	✓	✓	✓	✓	19	1	1	2	1	58.09
A2D2 [25]	2020	41,277	×	✓	✓	✓	✓	✓	✓	38	2	1	3	3	58.06
Virtual KITTI 2 [9]	2020	20,992	✓	✓	✓	✓	✓	✓	✓	8	4	2	3	1	56.74
RADIATE [26]	2021	44,140	✓	✓	✓	✓	✓	✓	✓	8	5	2	4	1	56.70
GTSDB [73]	2013	900	×	✓	✓	✓	✓	✓	✓	4	2	2	3	1	56.68
Caltech Ped [74]	2009	250K	×	✓	✓	✓	✓	✓	✓	1	1	1	1	2	56.64
SHIFT [10]	2022	2.5M	✓	✓	✓	✓	✓	✓	✓	23	5	5	3	-	55.79
CAOS [7]	2022	13K	×	✓	✓	✓	✓	✓	✓	13	3	2	3	-	55.74
STF [75]	2020	13,500	×	✓	✓	✓	✓	✓	✓	1	4	2	3	4	54.91
KITTI-360 [52]	2021	150K	✓	✓	✓	✓	✓	✓	✓	37	1	1	1	1	54.77
ACDC [76]	2021	4,006	×	✓	✓	✓	✓	✓	✓	19	4	2	3	1	54.62
VPNet [77]	2017	20K	×	✓	✓	✓	✓	✓	✓	17	3	2	1	1	54.37
ONCE [29]	2021	1M	×	✓	✓	×	×	×	×	5	3	2	5	1	54.09
V2X															
DAIR-V2X [27]	2021	71,254	✓	✓	✓	✓	✓	✓	✓	10	2	2	2	1	49.21
V2XSet [11]	2022	11,447	✓	✓	✓	✓	✓	✓	✓	1	1	1	1	-	45.39
V2VNet [78]	2020	51.2K	✓	✓	✓	✓	✓	✓	✓	1	1	1	1	-	44.94
T&J [79]	2019	100	✓	✓	✓	✓	✓	✓	✓	1	1	1	2	1	44.58
Rope3D [80]	2022	50K	×	✓	✓	✓	✓	✓	✓	12	3	3	2	1	44.39
V2X-Sim [81]	2022	10K	✓	✓	✓	✓	✓	✓	✓	23	1	1	1	-	41.00
V2V4Real [82]	2023	40K	✓	✓	✓	✓	✓	✓	✓	5	2	1	2	1	40.28
Co-Percep [83]	2020	10K	×	✓	✓	✓	✓	✓	✓	1	1	1	1	1	40.06
A9-Dataset [84]	2022	1,098	✓	✓	✓	✓	✓	✓	✓	9	4	2	1	1	39.12
LUMPI [85]	2022	200K	✓	✓	✓	✓	✓	✓	✓	3	6	3	1	1	35.87
DOLPHINS [86]	2022	42,376	✓	✓	✓	✓	✓	✓	✓	2	3	2	4	-	35.21
Drone															
UAVDT [87]	2018	80K	✓	✓	✓	✓	✓	✓	✓	3	2	2	6	1	61.63
DroneVehicle [88]	2021	28,439	×	✓	✓	✓	✓	✓	✓	5	1	3	4	1	44.42
Others															
Mapillary Vistas [89]	2017	25K	×	✓	✓	✓	✓	✓	✓	66	5	3	3	2	68.63
TT 100K [15]	2016	100K	×	✓	✓	✓	✓	✓	✓	45	2	2	2	10	61.99
Pascal3D+ [90]	2014	30,899	×	✓	✓	✓	✓	✓	✓	12	2	2	1	-	58.07
WildDash [91]	2018	1,800	×	✓	✓	✓	✓	✓	✓	28	2	2	7	1	50.02
TorontoCity [92]	2016	56K	×	✓	✓	✓	✓	✓	✓	4	2	2	1	1	45.82
DAWN [93]	2020	4,543	×	✓	✓	✓	✓	✓	✓	5	6	3	3	-	43.99
RAD [94]	2019	60	×	✓	✓	✓	✓	✓	✓	19	1	1	1	1	37.86
STCrowd [95]	2022	10,891	✓	✓	✓	✓	✓	✓	✓	1	3	1	1	1	35.19

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.

pixel-wise semantic labels, critical for training and validating autonomous vehicles' perception and navigation systems. ApolloScapes supports images and point clouds semantic segmentation, 2D/3D object detection, multi-object tracking, and lane segmentation, enabling the creation and evaluation of advanced and safe autonomous driving systems.

SemanticKITTI. SemanticKITTI [6] is a notable extension of the KITTI family, focusing on semantic segmentation in the autonomous driving field. SemanticKITTI consists of over 43,000 LiDAR point cloud frames, making it one of the largest 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 evaluating the performance of point cloud semantic segmentation methods, underpinning numerous studies and innovations in related domains.

nuScenes. nuScenes [22] stands as an essential contribution to the field of autonomous driving, offering a rich repository of data that addresses the diverse needs of perception systems.

nuScenes leverages LiDAR, radars, and cameras to record data from different urban scenes from Boston and Singapore. It is worth mentioning that its six cameras provide a comprehensive perspective of the surrounding environment, making them widely utilized in multi-view object detection tasks. In conclusion, the nuScenes dataset is a cornerstone in developing autonomous driving technologies, supporting multi-tasks and applications, and setting new benchmarks in the field.

Waymo. The Waymo Open Dataset [23], introduced in 2019, significantly influences research and advancement in autonomous driving by providing an extensive size of multimodal sensory data with high-quality annotations. 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, such as detection, tracking, and segmentation.

BDD100K. BDD100K [24] dataset, released by the Berkeley DeepDrive center in 2018, is a substantial and diverse driving dataset renowned for its size and diversity. It comprises

100,000 videos, each about 40 seconds in duration. Meanwhile, it offers various annotated labels for object detection, tracking, semantic segmentation, and lane detection. This extensive compilation of data prompts advancements in the autonomous driving community, establishing itself as a challenging and versatile platform for researchers and engineers to propose and refine algorithms.

RADIATE. RADIATE [26] is the first public radar dataset, which contains 44,140 annotated frames gathered under several adversarial weather conditions, such as rain, fog, overcast, and snow. It also incorporates LiDAR and camera data, allowing full perception and understanding of the driving surroundings.

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. In summary, the Argoverse 2 dataset offers numerous multimodal data of real-world driving scenarios, fostering innovations and advancements of algorithms and showing its substantial potential as a vital resource in autonomous driving.

2) **V2X: V2VNet.** The dataset introduced by V2VNet [78] focuses on leveraging V2V communication to allow autonomous vehicles to share information and perceive the environment from multiple viewpoints, which is crucial for detecting occluded objects and predicting the behavior of other traffic participants. This dataset was created using a high-fidelity LiDAR simulator named Lidarsim [96] that utilizes real-world data to generate realistic LiDAR point clouds for various traffic scenes. In conclusion, this work brings attention to the V2V as a promising avenue for improving the capabilities of autonomous vehicles.

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. The impact of the DAIR-V2X dataset on autonomous driving is significant because it sets a benchmark for complexities of vehicle-infrastructure cooperation thanks to its multiple scenarios from the real world.

Rope3D. Rope3D presented in [80] is an essential contribution to the perception systems, which addresses the critical gap in autonomous driving by leveraging data gathered from roadside cameras. Rope3D consists of 50,000 images under various environmental conditions, including different illumination (daytime, night, dusk) and weather scenarios (rainy, sunny, cloudy), ensuring a high diversity of data that reflects real-world complexities. Overall, the Rope3D dataset is a pioneering work that accelerates advancements in roadside perception for autonomous driving while being a critical tool for researchers and engineers to develop more robust and intelligent autonomous driving systems.

V2V4Real. V2V4Real [82] is the first large-scale, real-

world dataset to address V2V cooperative perception. The data was collected from two vehicles with multimodal sensors like LiDAR and cameras. V2V4Real focuses on a range of perception tasks, such as cooperation 3D object detection, cooperative 3D object tracking, and Sim2Real domain adaptation. This versatility makes it an invaluable resource for developing and benchmarking autonomous driving algorithms.

3) **Drone: UAVDT.** The UAVDT [87] 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 [88] 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+.** Pascal3D+ [90] is an extension of the PASCAL VOC 2022 [97], overcoming the limitations of previous datasets by providing a richer and more varied set of annotations for images. Pascal3D+ augments 12 rigid object categories, such as cars, buses, and bicycles, with 3D pose annotations and adds more images from ImageNet [98], resulting in a high degree of variability.

TT 100K. Tsinghua-Tencent 100K [15] solves the challenges of detecting and classifying traffic signs in realistic driving conditions. It provides 100K images, including 30,000 traffic-sign instances. Beyond its large data size, the high-resolution images encompass a wide range of illumination and weather conditions, making it robust for training and validating traffic sign recognition.

Mapillary Vistas. Mapillary Vistas dataset, proposed by [89] 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 training and evaluation of driving systems to forecast traffic dynamics, pedestrian movements, and other essential factors that influence driving decisions. By simulating myriad driving scenarios, they empower autonomous vehicles to make informed decisions, navigate complex environments, and maintain safety and efficiency on the road. Hence, we demonstrate in detail several high-impact datasets related to these tasks according to the data size, modalities, and citation number. We summarize the prediction, planning, and control datasets into task-specific and multi-task groups.

Dataset	Year	Sensing domain	Size	Tasks	Weather conditions	Time of day	Scenario conditions
Task-Specific							
JAAD [16]	2017	onboard	75K frames	pedestrian IT	sunny, rainy, cloudy, snowy	day, afternoon, night	urban
Dr(eye)ve [53]	2018	onboard	500K frames	driver's attention prediction	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 [99]	2019	onboard	24K trajectories	driver IT			5 intersections
Argoverse [19]	2019	onboard	300K trajectories	TP	variety	variety	urban
inD [100]	2020	drone	11.5K trajectories	road user prediction	sunny	day	4 urban intersections
PePscenes [101]	2020	onboard	719 frames	pedestrian BP			
openDD [102]	2020	drone	84,774 trajectories	pedestrian BP			7 roundabouts
nuPlan [103]	2021	drone	1.5K hours data	MPlan			4 cities
exiD [104]	2022	drone	16 hours data	TP	sunny	daytime	7 locations on highway
MONA [105]	2022	drone	702K trajectories	TP	sunny, overcast, rain	8 am to 5 pm	urban
Multi-Task							
INTERACTION [20]	2019	drone, V2X	110K trajectories	MPlan and MP, DM			(un)signalized intersection
BLVD [106]	2019	onboard	120K frames	4D OT, 5D event recognition		day, night	urban, highway
roundD [107]	2019	drone	13,746 road users	scenario classification, BP	sunny	daytime	(sub-)urban
Lyft Level 5 [108]	2021	drone	1.1K hours data	MPlan, MP			suburban
LOKI [109]	2021	onboard	644 scenarios	TP, BP			(sub-)urban
SceNDD [110]	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
Talk2BEV [111]	2023	onboard	20K QA pairs	DM, IT, spatial reasoning			
V2X-Seq (forecasting) [50]	2023	V2X	50K scenarios	online/offline VIC TP			28 urban intersections

TABLE III

PREDICTION, PLANNING, AND CONTROL DATASETS. WE DEMONSTRATE SEVERAL CRUCIAL DATASETS RELATED TO PREDICTION, PLANNING, AND CONTROL. BP: BEHAVIOR PREDICTION, IP: INTENTION PREDICTION, MP: MOTION PREDICTION, TP: TRAJECTORY PREDICTION, MPLAN: MOTION PLANNING, DM: DECISION-MAKING, OT: OBJECT TRACKING, QA: QUESTION-ANSWERING

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 110,000 cars and trucks. The dataset aims to overcome the limitations of existing measurement methods for scenario-based safety validation, which often fail to capture the naturalistic behavior of road users or to include all relevant data with sufficient quality.

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. The PIE dataset offers rich annotations for perception and visual reasoning, including bounding boxes with occlusion flags, crossing intention confidence, and text labels for pedestrian actions. The long continuous sequences and annotations facilitate multiple tasks like trajectory prediction and pedestrian intention prediction.

USyd. USyd [99] pushes the progress of the driver intention prediction in the context of urban intersections without traffic signals, which are common in urban settings and represent a challenge due to the lack of clear road rules and signals. The dataset incorporates over 23,000 vehicles traversing five different intersections, collected using a vehicle-mounted LiDAR-based tracking system. The data modalities included are exhaustive, providing vehicle tracks with lateral and longitudinal coordinates, heading, and velocity. This information is vital to predict driver behavior, accounting for the uncertainty inherent in human driving patterns.

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.

inD. The importance of the inD [100] lies in its large-scale, high-quality, and diverse set of trajectory data crucial for several applications, including road user prediction models and scenarios-based safety validation for autonomous vehicles in urban intersection environments. It covers around 11,500 different road user trajectories, e.g., vehicles, bicyclists, and pedestrians. These trajectories have a positioning error of less than 0.1 meters, which is pivotal for the reliability of the data.

PePscenes. PePscenes [101] tackles the requirement for understanding and anticipating pedestrian actions in dynamic driving surroundings. This dataset enhances the nuScenes [22] dataset by adding per-frame 2D/3D bounding box and behavior annotations, focusing on pedestrian crossing actions. One of the key attributes of [101] is incorporating various data types involving semantic maps, scene images, trajectories, and ego-vehicle states, which are necessary for creating robust models capable of understanding complex traffic scenarios.

openDD. The openDD [102] dataset focuses on analyzing and predicting traffic scenes around roundabouts, which are complex and unregulated by traffic lights. It is unique in leveraging drone-captured imagery with high resolution (4K), which spans over 62 hours of trajectory data from 501 separate flights. The dataset contains not only trajectories but also shapefiles and an extensible markup language (XML) file describing the road topology, along with a reference image for each underlying intersection.

nuPlan. nuPlan [103] 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, featuring 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.

exiD. The exiD [104] 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 in terms of the diversity of interactions captured, especially those involving lane changes at highway entries and exits.

MONA. The Munich Motion Dataset of Natural Driving (MONA) [105] 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.

2) *Multi-Task Datasets:* **INTERACTION.** The INTERACTION [20] dataset covers diverse, complex, and critical driving scenes, coupled with its comprehensive semantic map, allowing it to be a versatile platform for a multitude of tasks, such as motion prediction, imitation learning, and validation of decision and planning. Its inclusion of different countries further improves the robustness of analyzing the driving behavior across different cultures, which is critical for the global development of autonomous driving.

BLVD. The BLVD [106] benchmark facilitates tasks such as dynamic 4D (3D+temporal) tracking, 5D (4D+interactive) interactive event recognition, and intention prediction, which are essential for a deeper understanding of traffic scenes. BLVD offers around 120K frames from different traffic scenes, comprising object density (low and high) and illumination conditions (daytime and nighttime). These frames are fully annotated with a large number of 3D labels incorporating vehicles, pedestrians, and riders.

round. The round dataset presented by [107] is pivotal for scenario classification, road user behavior prediction, and driver modeling because of the large number of collections 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.

Lyft Level 5. Lyft Level 5 [108] 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 the Lyft Level 5 dataset a vital benchmark for prediction and planning.

LOKI. LOKI [109], 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.

SceNDD. SceNDD [110] introduces real driving scenarios that feature diverse trajectories and driving behaviors for developing efficient motion planning and path-following algorithms. It is also adaptable to different configurations of the ego-vehicle and contains a time horizon of prediction that can be broken down into timestamps for detailed analysis. In conclusion, the SceNDD dataset is a significant addition to autonomous driving prediction and planning research.

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 for improving the predictive capabilities of autonomous driving systems 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 bird's-eye-view (BEV) semantic segmentation.

Talk2Bev. The innovative dataset, Talk2BEV [111], promotes the shift from traditional autonomous driving tasks to incorporate large vision-language models with BEV maps in the context of autonomous driving. Talk2BEV utilizes recent advances in vision-language models, allowing for a more flexible and comprehensive understanding of road scenarios. This dataset comprises over 20,000 diverse question categories, all human-annotated and derived from [22]. The proposed Talk2BEV-Bench benchmark can be leveraged across multiple tasks, covering decision-making, visual and spatial reasoning, and intent prediction.

V2X-Seq (Forecasting). The trajectory forecasting dataset is a substantial part of the real-world dataset V2X-Seq [50], containing about 80K infrastructure-view and 80K vehicle-view scenarios, and a further 50K cooperative-view scenarios. This diversity of sensing domains creates a more holistic view of the traffic environment, exhibiting huge potentiality for research and analysis on vehicle-infrastructure-cooperative (VIC) trajectory prediction.

C. End-to-End Datasets

End-to-end has become a growing trend as an alternative to modular-based architecture in autonomous driving [66]. Several versatile datasets (like nuScenes [22] and Waymo [23]) or simulators like CARLA [112] 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 [113] dataset is notable for its use of event-based cameras, which provide 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.

Other datasets summarized in this survey are shown in Tab. IV Tab. V, Tab. VI.

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. How annotations are created

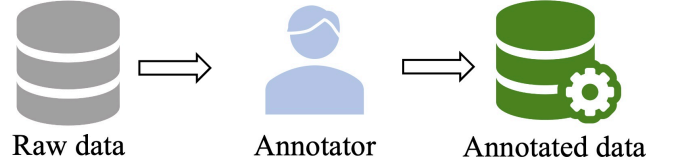
Different AD tasks require specific types of annotation. For example, object detection requires bounding box labels of instance, segmentation is based on pixel- or point-level annotations, and continually labeled trajectory is critical for trajectory prediction. On the other hand, as shown in Fig. 11,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

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.

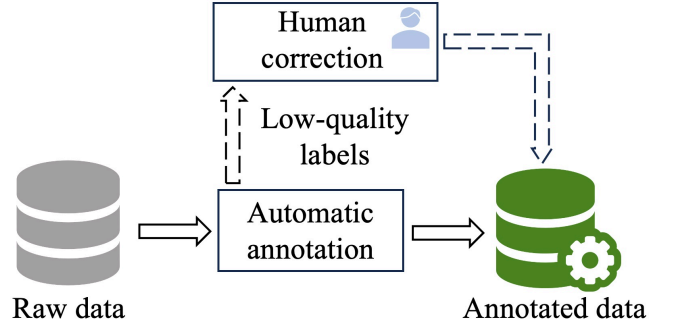
Many studies have proposed fully or semi-automatic annotation methods to improve annotation efficiency. [264] presented a fully automatic annotation approach based on weakly supervised learning to segment proposed drivable paths in images. [265] is a semi-automatic annotation method utilizing the objectness priors to generate segmentation masks. After that, [266] offered a semi-automatic method considering 20 classes. Polygon-RNN++ [267] presented an interactive segmentation annotation tool following the idea of [268]. Instead of using image information to generate pixel-level labels, [269] explored transferring 3D information into 2D image domains to generate semantic segmentation annotations. For labeling 3D data, [270] proposed an image-assisted annotation pipeline. [271] leveraged active learning to select a few points and to form a minimal training set to avoid labeling the whole point cloud scenes. [272] introduced an efficient labeling framework with semi/weakly supervised learning to label outdoor point clouds.

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.

(a) Manual annotation



(b) Semi-automatic annotation



(c) Fully automatic annotation

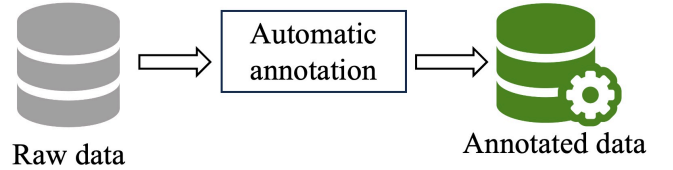


Fig. 11. Annotation pipelines. We demonstrate (a) Manual annotation: The professional annotators label the raw data using annotation tools. (b) Semi-automatic annotation: After generating annotations using an automatic annotation algorithm, the low-quality labels are refined by annotators. (c) Fully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

Labelme [273] 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. [274] presented a semi-automatic video labeling tool based on the open-source video annotation system VATIC⁷. [275] is another automatic video annotation tool for AD scenes. Compared to daytime annotation, process bounding box annotations in the nighttime is more challenging. [276] 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. [277] applied a two-stage weakly supervised learning framework using human-in-the-loop to label LiDAR point clouds. ViT-WSS3D [278] generated pseudo-bounding boxes by modeling global interactions between LiDAR points and corresponding weak labels. Apolloscape [17] employed a labeling pipeline similar to [269], which consists of a 3D labeling and a 2D labeling branch, to handle static background/objects and moving objects, respectively. 3D BAT [279] developed an annotation

⁷<https://www.cs.columbia.edu/~vondrick/vatic/>

Dataset	Year	Size	Temporal	Sensing domain	Tasks	Real/Synthetic
KAIST MPD [114]	2015	95K color-thermal pair frames	×	onboard	pedestrian detection	real
BelgiumTS [115]	2011	13K traffic sign annotations	×	onboard	traffic sign detection	real
nighttime drive [116]	2018	35K	×	onboard	SS	real
D ² -City [117]	2019	700K annotated frames	✓	onboard	OD, MOT	real
Caltech Lanes [118]	2019	1,224 frames	×	onboard	lane detection	real
A*3D [119]	2020	39,179 point cloud frames	×	onboard	(3D) OD	real
PreSIL [8]	2019	50K frames	✓	onboard	OD, 3D SS	synthetic
H3D [120]	2019	27,721 frames	✓	onboard	(3D) OD, MOT	real
ROAD [121]	2021	122K frames	✓	onboard	OD, SS	real
ETH Ped [122]	2007	2,293 frames	✓	onboard	pedestrian detection	real
CADC [123]	2020	56K	×	onboard	(3D) OD, OT	real
RadarScenes [124]	2021	40K radar frames	✓	onboard	OD, classification	real
All-in-One Drive [125]	2021	100K	✓	onboard	(3D) OD, (3D) SS, trajectory prediction	real
CURE-TSR [126]	2017	2.2M annotated images	×	onboard	traffic sign detection	real
Paris-Lille-3D [127]	2018	2,479 frames	×	onboard	3D SS, classification	real
PandaSet [128]	2021	8,240 frames	✓	onboard	(3D) OD, SS, OT	real
NightOwls [129]	2018	56K frames	✓	onboard	pedestrian detection, tracking	real
SynWoodScape [130]	2022	80K frames	×	onboard	(3D) OD, segmentation	synthetic
Boreas [131]	2023	7,111 frames	✓	onboard	(3D) OD, localization	real
TUD-Brussels [132]	2009	1,600 frames	×	onboard	(3D) OD	real
VEIS [133]	2018	61,305 frames	×	onboard	OD, SS	synthetic
CCTSDb 2021 [134]	2021	16,356 frames	×	onboard	traffic sign detection	real
SemanticPOSS [135]	2020	2,988 point cloud frames	✓	onboard	3D SS	synthetic
IDDA [136]	2020	1M frames	×	onboard	segmentation	synthetic
TJ4RadSet [137]	2022	7,757 frames	✓	onboard	OD, OT	real
TME Motorway [138]	2012	30K frames	✓	onboard	OD, OT	real
Stanford Tack [139]	2011	14K tracks	✓	onboard	classification	real
CARRADA [140]	2020	7,193 radar frames	✓	onboard	SS	real
SODA10M [141]	2021	20K labeled images	✓	onboard	OD	real
Robo3D [142]	2023	476K frames	✓	onboard	(3D) OD, 3D SS	real
LostAndFound [143]	2016	2,104 frames	✓	onboard	obstacle detection	real
Titan [144]	2020	75,262 frames	✓	onboard	OD, action recognition	real
CODA [145]	2022	1,500 frames	×	onboard	corner case detection	real
PixSet [146]	2021	29K point cloud frames	×	onboard	(3D) OD	real
ZOD [30]	2023	100K frames	✓	onboard	(3D) OD, segmentation	real
K-Radar [147]	2023	35K radar frames	×	onboard	3D OD, OT	real
RoadObstacle21 [148]	2021	327 scenes	×	onboard	anomaly segmentation	synthetic
VIL-100 [149]	2021	10K frames	×	onboard	lane detection	real
EuroCity Persons [150]	2018	47,300 frames	×	onboard	OD	real
TuSimple [151]	2017	6,408 frames	×	onboard	lane detection, velocity estimation	real
OpenMPD [152]	2021	15K frames	×	onboard	(3D) OD, 3D OT, semantic segmentation	real
WADS [153]	2021	1K point cloud frames	✓	onboard	SS	real
NightCity [154]	2020	4,297 frames	×	onboard	nighttime SS	real
LiDAR Snowfall [155]	2022	7,385 point cloud frames	✓	onboard	3D OD	synthetic
LLMAS [156]	2019	100K frames	✓	onboard	lane detection	real
Fishyscapes [157]	2021	1,030 frames	×	onboard	SS, anomaly detection	real
Toronto-3D [158]	2020	4 scenarios	×	onboard	3D SS	real
MIT-AVT [159]	2020	1.15M 10s video clips	✓	onboard	SS, anomaly detection	real
LISA [160]	2014	6,610 frames	✓	onboard	traffic sign detection	real
SynthCity [161]	2019	367.9M points in 30 scans	×	onboard	(3D) OD, (3D) SS	synthetic
SemanticUSL [162]	2021	1.2K frames	×	onboard	domain adaptation 3D SS	real
RANUS [163]	2018	4K frames	×	onboard	SS, scene understanding	real
PePScenes [101]	2020	40K frames	✓	onboard	(3D) OD, pedestrian action prediction	real
MUAD [164]	2022	10.4K frames	✓	onboard	OD, SS, depth estimation	synthetic
AUTOCASTSIM [165]	2022	52K frames	✓	V2X	(3D) OD, OT, SS	real
Mcity [166]	2019	1,7500 frames	✓	onboard	SS	real
comma2k19 [167]	2018	2M images	×	onboard	pose estimation, end-to-end driving	real
DET [168]	2019	5,424 event-based camera images	×	onboard	lane detection	real
TRoM [169]	2017	712 frames	×	onboard	road marking detection	real
Cirrus [170]	2020	6,285 frames	✓	onboard	(3D) OD	real
KITTI InstanceMotSeg [171]	2020	12,919 frames	✓	onboard	moving instance segmentation	real
CARTI [172]	2022	11K frames	✓	V2X	cooperative perception	synthetic
Multifog KITTI [173]	2021	15K frames	×	onboard	3D OD	synthetic
UNDD [174]	2019	7.2K frames	✓	onboard	SS	real
CeyMo [175]	2021	2,887 frames	×	onboard	road marking detection	real
Raidar [176]	2021	58,542 rainy street scenes	×	onboard	SS	real
K-Lane [177]	2022	15,382 frames	×	onboard	lane detection	real
aiMotive [178]	2023	26,583 frames	✓	onboard	3D OD, MOT	real
SAP [179]	2016	19K frames	×	drone	OD, OT	real
Astyx [180]	2019	500 radar frames	×	onboard	3D OD	real
Comap [181]	2021	4,391 frames	✓	V2X	3D OD	synthetic
Ithaca365 [182]	2022	7K frames	×	onboard	3D OD, SS, depth estimation	real
Small Obstacles [183]	2020	3K frames	×	onboard	small obstacle segmentation	real
GLARE [184]	2022	2,157 frames	×	onboard	traffic sign detection	real
WIBAM [185]	2021	33,092 frames	✓	V2X	3D OD	real
MIT DriveSeg [186]	2019	5K frames	×	onboard	SS	real
SUPS [187]	2023	5K frames	✓	onboard	SS, depth estimation, SLAM	synthetic
R3 [188]	2021	369 scenes	×	onboard	out-of-distribution detection	real

TABLE IV
AUTONOMOUS DRIVING DATASET I

Dataset	Year	Size	Temporal	Sensing domain	Tasks	Real/Synthetic
V2X-Seq (perception) [50]	2023	15K frames	✓	V2X	cooperative perception	real
DRIV100 [189]	2021	100K frames	×	onboard	domain adaptation SS	real
NEOLIX [190]	2021	30K frames	✓	onboard	3D OD, OT	real
CARLANE [191]	2022	118K frames	✓	onboard	lane detection	synthetic
IPS3000+ [192]	2021	14,198 frames	✓	V2X	3D OD	real
TUM-Traffic [193]	2023	4.8K frames	✓	V2X	3D OD, OT	real
Amodal Cityscapes [194]	2022	5K frames	×	onboard	amodal SS	real
WEDGE [195]	2023	3,360 frames	×	others	OD, classification	synthetic
HEV [196]	2019	230 video clips	✓	onboard	object localization	real
Aachen Day-Night [197]	2018	4,328 images, 1.65M points	✓	onboard	visual localization	real
HAD [198]	2019	5,675 video clips	✓	onboard	end-to-end driving	real
CARLA-100 [199]	2019	100 hours driving	✓	onboard	path planning, behavior cloning	synthetic
Collective activity [200]	2009	44 short videos	✓	others	human activity recognition	real
DA4AD [201]	2020	9 sequences	✓	onboard	visual localization	real
Bosch STL [202]	2017	13,334 images	✓	onboard	traffic light detection and classification	real
UrbanLaneGraph [203]	2023	around 5,220 km lane spans	✓	drone	lane graph estimation	real
CrashD [204]	2022	15,340 scenes	×	onboard	3D OD	synthetic
CODD [205]	2022	108 sequences	✓	V2X	multi-agent SLAM	synthetic
MIT DSS [186]	2020	10K video frames	✓	onboard	SS	real
CPIS [83]	2020	10K frames	×	V2X	cooperative 3D OD	synthetic
CarlaScenes [206]	2022	7 sequences	✓	onboard	(3D) SS, SLAM, depth estimation	synthetic
FlyingThings3D [207]	2015	26,066 frames	✓	others	scene flow estimation	synthetic
San Francisco Landmark [208]	2011	150K panoramic images	×	others	landmark identification	real
DDAD [209]	2020	21,200 frames	×	onboard	depth estimation	real
Ua-detrac [210]	2015	140K frames	✓	V2X	OD, MOT	real
NCLT [211]	2015	34.9 hours	✓	onboard	odometry	real
NVSEC [212]	2018	around 28 km distances	✓	others	SLAM	real
MSLU [213]	2013	36.8 km distances	✓	onboard	SLAM	real
Oxford Radar RobotCar [214]	2020	240K scans	✓	onboard	odometry	real
A2D2 [25]	2020	433,833 frames	✓	onboard	3D OD, SS, SLAM	real
VERI-Wild [215]	2019	416,314 images	✓	V2X	onboard re-identification	real
OTOH [108]	2020	170K scenes	✓	drone	trajectory prediction, planning	real
TRANCOS [216]	2015	1.2K images	×	V2X	onboard number estimation	real
Complex Urban [217]	2017	around 190km paths	✓	onboard	SLAM	real
Syncapes [218]	2018	25K images	×	onboard	OD, SS	synthetic
SydneyUrbanObject [219]	2013	588 object scans	×	onboard	classification	real
ApolloCar3D [220]	2019	5,277 driving images	×	onboard	3D instance understanding	real
ACDC [76]	2021	4,006 images	×	onboard	SS on adverse conditions	real
MulRan [221]	2020	41.2km paths	✓	onboard	place recognition	real
Paris-rue-Madame [222]	2014	643 objects	×	onboard	OD, SS	real
Ground Truth SitXel [223]	2013	78,500 frames	✓	onboard	stereo confidence	real
LiVi-Set [224]	2018	10K frames	✓	onboard	driving behavior prediction	real
Newer College [225]	2020	290M points, 2300 seconds	✓	others	SLAM	real
CADP [226]	2018	1,416 scenes	×	V2X	traffic accident analysis	real
LIBRE [227]	2020	40 frames	×	onboard	LiDAR performance benchmark	real
OPV2V [228]	2022	11,464 frames	×	V2X	onboard-to-onboard perception	synthetic
NYC3DCars [229]	2013	2K images	×	onboard	OD	real
RUGD [230]	2019	37K frames	✓	others	SS	real
EU LTD [231]	2020	around 37 hours	✓	onboard	odometry	real
MOTSynth [232]	2021	768 driving sequences	✓	onboard	Pedestrian detection and tracking	synthetic
LLamas [156]	2019	100,042 images	×	onboard	lane marker detection, lane segmentation	real
PedX [233]	2018	5K images	✓	onboard	pedestrian detection and tracking	real
CCD [234]	2020	4.5K videos	✓	onboard	accident prediction	real
MAVD [235]	2021	113,283 images	✓	onboard	OD and OT with sound	real
Gated2Depth [236]	2020	17,686 frames	✓	onboard	depth estimation	real
DDD20 [237]	2017	51 hours event frames	✓	onboard	end-to-end driving	real
TCGR [238]	2020	839,350 frames	✓	others	traffic control gesture recognition	real
4Seasons [238]	2020	350km recordings	✓	onboard	SLAM	real
CCSAD [239]	2015	96K frames	✓	onboard	scene understanding	real
TAF-BW [240]	2018	2 scenarios	✓	V2X	MOT, V2X communication	real
Boxy [241]	2019	200K frames	✓	onboard	OD	real
AMUSE [242]	2013	117,440 frames	✓	onboard	SLAM	real
Brno Urban [243]	2019	375.7 km	✓	onboard	recognition	real
AUTOMATUM [244]	2021	30 hours	✓	drone	trajectory prediction	real
DurLAR [245]	2021	100K frames	✓	onboard	depth estimation	real
Reasonable-Crowd [246]	2021	92 scenarios	✓	onboard	driving behavior prediction	synthetic
Daimler Ped [247]	2013	12,485 frames	✓	onboard	pedestrian path prediction	real
PVDN [248]	2021	59,746 frames	✓	onboard	nighttime OD, OT	real
CARLA-WildLife [249]	2022	26 videos	✓	onboard	out-of-distribution tracking	synthetic
SOS [249]	2022	20 videos	✓	onboard	out-of-distribution tracking	real
RoadSaW [250]	2022	720K scenes	✓	onboard	Road surface and wetness estimation	real
I see you [251]	2022	170 sequences, 340 trajectories	✓	V2X	OD	real
ASAP [252]	2022	1.2M images	✓	onboard	online 3D OD	real
OpenLane-V2 [253]	2023	466K images	✓	onboard	lane detection, scene understanding	real
Daimler Stereo Ped [254]	2011	21,790 frames	✓	onboard	pedestrian detection	real
Brain4Cars [255]	2015	2M frames	✓	onboard	Maneuver Anticipation	real
NEXET [256]	2017	91,190 frames	×	onboard	OD	real
DIML [257]	2017	470 videos	✓	onboard	lane detection	real

TABLE V
AUTONOMOUS DRIVING DATASET 2

Dataset	Year	Size	Temporal	Sensing domain	Tasks	Real/Synthetic
UAH-Driveset [258]	2016	500 mins	✓	onboard	lane detection, detection	real
SSCBENCH [259]	2023	66,913 frames	×	onboard	semantic scene completion	real
FLIR [260]	-	26,442 thermal frames	✓	onboard	thermal image OD	real
VLMV [261]	2020	900 frames	✓	V2X	lane merge	real
CityFlow [262]	2019	200K bounding boxes	✓	V2X	OD, MOT, re-identification	real
PREVENTION [263]	2019	356 mins	✓	onboard	onboard trajectory and intention prediction	real

TABLE VI
AUTONOMOUS DRIVING DATASET 3

toolbox to assist in obtaining 2D and 3D labels in semi-automatic labeling.

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, [280] online generated actions for maneuvers and were annotated into the trajectory. [281] consists of a crowd-sourcing step followed by a precise process of expert aggregation. [282] developed an active learning framework to annotate driving trajectory. Precisely anticipating movement patterns of pedestrians is critical for driving safety. Styles et al. [283] introduced 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. To achieve this goal, some researchers utilized the Grand Theft Auto 5 (GTA5) game engine to build datasets [4], [5]. [284] presented a real-time system based on multiple games to generate annotations for various AD tasks. Instead of applying game videos, SHIFT [10], CAOS [7], and V2XSet [11] were created based on the CARLA [112] simulator. Compared to [11], V2X-Sim [81] studied employing multiple simulators [112], [285] to generate dataset for V2X perception tasks. CODD [205] further exploited using [112] to generate 3D LiDAR point clouds for cooperative driving. Other works [2], [3], [9], [133] leveraged the Unity development platform [286] to generate synthetic datasets.

B. The quality of annotations

Existing supervised learning-based AD algorithms build upon numerous labeled data. However, training on low-quality annotations can negatively affect the safety and reliability of autonomous vehicles. Therefore, ensuring the quality of annotations is fundamental for improving accuracy while driving in complex real-world environments. According to

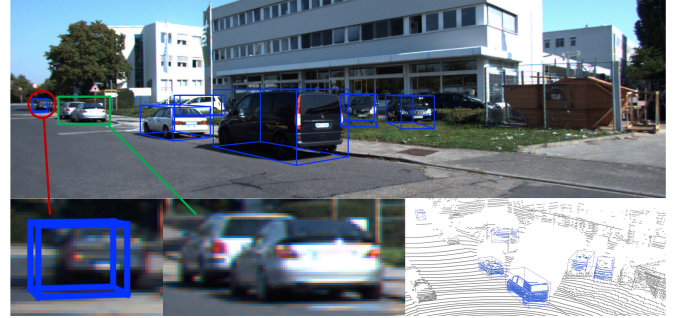


Fig. 12. Mislabeling example of KITTI [13] dataset. We show the ground truth in blue. The bounding box of a car (in red circle) is not precise. Two cars (in green cube) are not annotated, although sensors obviously capture them.

the study [287], 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 dataset purpose 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. [288] 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 cuboid) are not annotated, even though the camera and LiDAR capture them clearly.

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

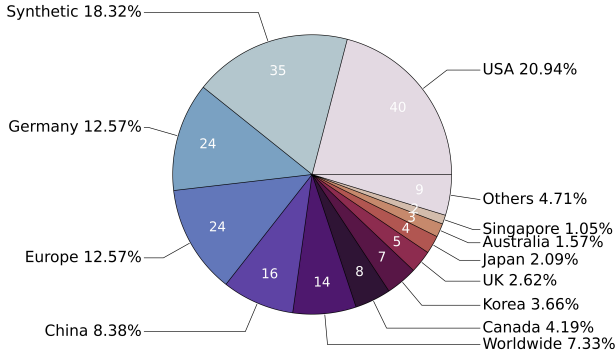


Fig. 13. The distribution of datasets around the world. This figure illustrates the distribution of data collection locations of the datasets.

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, reflecting its robust automotive industry and influence on the advancement of autonomous technology. China is closely followed with 16 datasets, suggesting the interest and investment in this field of China. Another notable point is that there are 11 datasets collected worldwide and 24 from the European region (except Germany). This diverse regional distribution enhances the robustness of the collected data and highlights the international efforts and collaborations in the research community and industry.

On the other hand, although smaller segments represent a range of other countries, including Canada, Korea, the UK, Japan, and Singapore, these countries are developed countries with solid technological backgrounds and accumulation—an extreme regional bias reflected by this statistic. The dominance of the USA, western Europe, and East Asia leads to the bias where autonomous driving systems are overfitted to environmental conditions typical of these regions. This bias could result in autonomous vehicles’ failure to perform under varied or unseen regions and cases. Hence, introducing data from a more expansive array of countries and regions, such as Africa, can assist in the comprehensive development of autonomous vehicles.

Furthermore, 35 synthetic datasets generated by simulators like CARLA [112] 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.

B. Chronological Trends in Perception Datasets

In Fig. 10, we introduce a chronological overview of perception datasets with the top 50 impact scores from 2007 to 2023

(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] (2021) also exhibit a trend toward cooperative driving systems. Furthermore, because of the non-occlusion perspective provided by UAV, drone-based datasets, such as UAVDT [87] published in 2018, take a crucial position in advancing perception systems.

C. Data Distribution

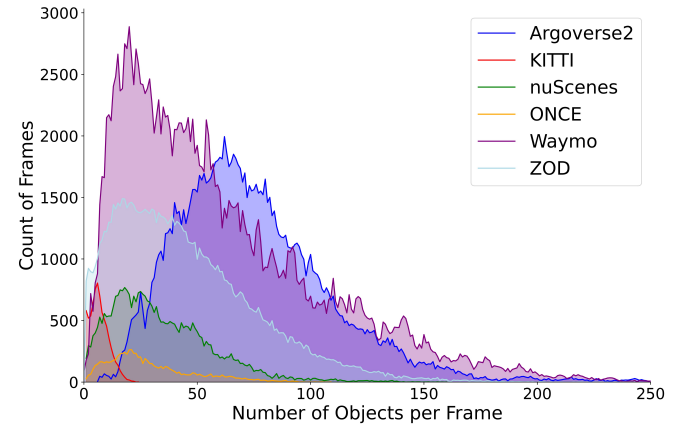


Fig. 14. Comparison of the distribution of the number of objects per frame across several datasets: Argoverse 2 [28], KITTI [13], nuScenes [22], ONCE [29], Waymo [23], and ZOD [30]. The horizontal axis quantifies the number of objects detected in a single frame, while the vertical axis represents the count of frames containing that number of objects.

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 [29], its density of objects evenly distributes in the supported perception range. Datasets like nuScenes [22] and ZOD [30] demonstrate similar

Dataset	Year	Size	Temporal	Sensing domain	Tasks	Real/Synthetic
BDD-X [289]	2018	8.4M	✓	onboard	reasoning, planning	real
Cityscapes-Ref [290]	2018	5,000 stereo videos	✓	onboard	object referring	real
TOUCHDOWN [291]	2019	9,326 examples	✓	onboard	reasoning, navigation	real
Talk2Car [292]	2019	11,959 commands, 850 videos	✓	onboard	object referring	real
BDD-OIA [293]	2020	11,303 scenarios	×	onboard	explainable decision-making	real
CityFlow-NL [294]	2021	5,289 samples	✓	onboard	OT	real
CARLA-NAV [295]	2022	83K	✓	onboard	navigation	real
NuPrompt [296]	2023	34K frames, 35K prompts	✓	onboard	MOT	real
NuScenes-QA [297]	2023	34K scenes, 460K QA pairs	✓	onboard	visual question answering	real
Refer-KITTI [51]	2023	6,650 frames	✓	onboard	referring MOT	real
Driving LLMs [298]	2023	10K driving situations, 160K QA pairs	×	drone	visual question answering	synthetic
DRAMA [299]	2023	17,785 scenarios	✓	onboard	reasoning, visual question answering	real
Rank2Tell [300]	2023	116 scenarios	✓	onboard	importance level ranking	real
LamPilot [301]	2023	4,900 samples	✓	others	planning	synthetic
LangAuto CARLAR [302]	2023	64K data clips	✓	onboard	closed-loop driving	synthetic
NuScenes-MQA [303]	2023	34K scenarios, 1.4M QA pairs	×	onboard	visual question answering	real
DriveMLM [304]	2023	280 hours	✓	onboard	planning, control	synthetic
DriveLM-nuScenes [305]	2023	4,871 frames	✓	onboard	end-to-end driving	real
DriveCARLA [305]	2023	183,373 frames	✓	onboard	end-to-end driving	real
LiDAR-text [306]	2023	420K 3D captioning data, 280K 3D grounding data	×	onboard	3D scene understanding	real

TABLE VII

VLM AUTONOMOUS DRIVING DATASET. OT: OBJECT TRACKING, MOT: MULTI-OBJECT TRACKING, QA: QUESTION-ANSWERING

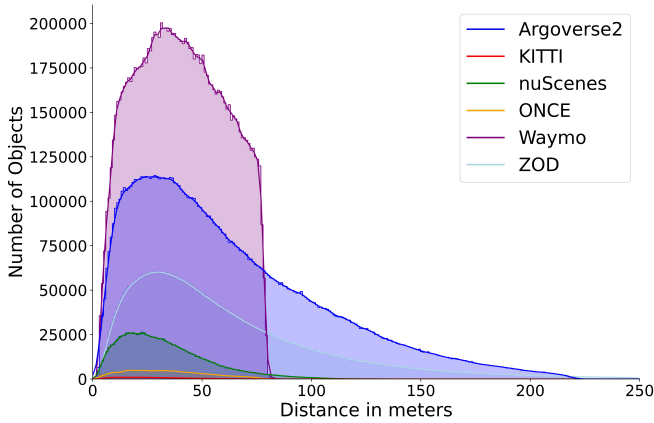


Fig. 15. Comparison of the distribution of the number of objects detected at various distances across several datasets: Argoverse 2 [28], KITTI [13], nuScenes [22], ONCE [29], Waymo [23], and ZOD [30]. The horizontal axis measures the distance from the ego vehicle in meters, and the vertical axis quantifies the number of objects detected at that distance.

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 tapers 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 detection.

VIII. DISCUSSION AND FUTURE WORKS

In this paper, we mainly focus on analyzing the existing datasets, which usually include rich visual data, and aim at tasks in the modular-based pipeline. However, with the rapid technological development, especially the excellent performance of the Large Language Models, many novel trends of the 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 the UniAD [307] verifies the potential ability of end-to-end models. However, the number of datasets for end-to-end AD is limited [113], [167]. 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 [67].

Introduce Language into AD Datasets. The 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. [308] highlighted the prominent role of Multimodal Large Language Models in various AD tasks, such as perception [290], [296], motion planning [301], and motion control [111]. The autonomous driving datasets including language labels are shown in Tab. VII. Overall, incorporating language into AD datasets is a trend of the future development of AD datasets.

Data Generation via VLMs. As mentioned by [309], the powerful capability of VLMs can be used to generate autonomous driving data. For example, DriveGAN [310] 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 [311]–[313] explored world models to generate high-quality driving videos. DriveDreamer [312] as a pioneering work derived from real-world scenarios, addressing the limitation of the gaming environments or simulated settings.

Domain Adaptation. Domain adaptation is a critical challenge in developing autonomous vehicles [314], 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 environmental conditions [315], sensor settings [316], or synthetic-to-real transform [317].

IX. CONCLUSION

In this paper, we exhaustively and systematically reviewed and analyzed more than 200 existing autonomous driving datasets. We started with the sensor types and modalities, sensing domains, and tasks relevant to autonomous driving datasets. We introduced a novel evaluation metric called impact score to validate the influence and importance of perception datasets. Afterward, we demonstrated several high-impact datasets involving perception, prediction, planning, control, and end-to-end autonomous driving. Additionally, we explained the annotation methodology for autonomous driving datasets and investigated the factors affecting annotation quality.

Moreover, we delineated the chronological and geographical distribution of the collected datasets, providing a comprehensive perspective for understanding the current development of autonomous driving datasets. Meanwhile, we studied the data distribution of several datasets, offering a specific viewpoint for comprehending the variances across different datasets. In the end, we discussed the development and trend of the next generation of autonomous driving datasets.

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. Kloecker, 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, “nusenes: 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. Kretschmar, X. Dotiwalla, A. Choudhury, 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. Muehlegg, 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] 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.
 - [30] 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.
 - [31] 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.
 - [32] 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.
 - [33] 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.
 - [34] 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.
 - [35] 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.
 - [36] 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.
 - [37] 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.
 - [38] 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.
 - [39] 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.
 - [40] 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.
 - [41] 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.
 - [42] B. Kitchenham, “Procedures for performing systematic reviews,” *Keele, UK, Keele University*, vol. 33, no. 2004, pp. 1–26, 2004.
 - [43] 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.
 - [44] 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.
 - [45] 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.
 - [46] G. Chen, H. Cao, J. Conradt, H. Tang, F. Rohrbach, 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.
 - [47] R. Gade and T. B. Moeslund, “Thermal cameras and applications: a survey,” *Machine vision and applications*, vol. 25, pp. 245–262, 2014.
 - [48] 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.
 - [49] 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.
 - [50] 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.
 - [51] 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.
 - [52] 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.
 - [53] 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.
 - [54] 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.
 - [55] 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.
 - [56] 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.
 - [57] 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.
 - [58] 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.
 - [59] J. Lambert and J. Hays, “Trust, but verify: Cross-modality fusion for hd map change detection,” *arXiv preprint arXiv:2212.07312*, 2022.
 - [60] 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.
 - [61] 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.
 - [62] 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.
 - [63] 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.
 - [64] N. Sharma, C. Dhiman, and S. Indu, “Pedestrian intention prediction for autonomous vehicles: A comprehensive survey,” *Neurocomputing*, 2022.
 - [65] 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.
 - [66] 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.
 - [67] 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.
 - [68] 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.
 - [69] 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.
 - [70] 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.
 - [71] G. Varma, A. Subramanian, A. Nambodiri, 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.
- [72] 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.
 - [73] 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.
 - [74] 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.
 - [75] M. Bjelic, 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.
 - [76] 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.
 - [77] S. Lee, J. Kim, J. Shin Yoon, S. Shin, O. Bailo, N. Kim, T.-H. Lee, H. Seok Hong, S.-H. Han, and I. So Kweon, “Vpnet: Vanishing point guided network for lane and road marking detection and recognition,” in *Proceedings of the IEEE international conference on computer vision*, pp. 1947–1955, 2017.
 - [78] 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.
 - [79] 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.
 - [80] 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.
 - [81] 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.
 - [82] 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.
 - [83] 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.
 - [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] 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.
 - [86] R. Mao, J. Guo, Y. Jia, Y. Sun, S. Zhou, and Z. Niu, “Dolphins: Dataset for collaborative perception enabled harmonious and interconnected self-driving,” in *Proceedings of the Asian Conference on Computer Vision*, pp. 4361–4377, 2022.
 - [87] 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.
 - [88] 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.
 - [89] G. Neuhold, T. Ollmann, S. Rota Bulò, 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.
 - [90] 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.
 - [91] O. Zendel, K. Honauer, M. Murschitz, D. Steininger, and G. F. Dominguez, “Willdash-creating hazard-aware benchmarks,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 402–416, 2018.
 - [92] 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.
 - [93] M. A. Kenk and M. Hassaballah, “Dawn: vehicle detection in adverse weather nature dataset,” *arXiv preprint arXiv:2008.05402*, 2020.
 - [94] 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.
 - [95] 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.
 - [96] S. Manivasagam, S. Wang, K. Wong, W. Zeng, M. Sazanovich, S. Tan, B. Yang, W.-C. Ma, and R. Urtasun, “Lidarsim: Realistic lidar simulation by leveraging the real world. 2020 ieee,” in *CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 11164–11173, 2020.
 - [97] 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.
 - [98] 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.
 - [99] 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.
 - [100] 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.
 - [101] A. Rasouli, T. Yau, P. Lakner, S. Malekmohammadi, M. Rohani, and J. Luo, “Pepsccenes: A novel dataset and baseline for pedestrian action prediction in 3d,” *arXiv preprint arXiv:2012.07773*, 2020.
 - [102] A. Breuer, J.-A. Termöhlen, S. Homocanu, 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.
 - [103] 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.
 - [104] 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.
 - [105] 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.
 - [106] 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.
 - [107] 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.
 - [108] 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.
 - [109] 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.
 - [110] A. Prabu, N. Ranjan, L. Li, R. Tian, S. Chien, Y. Chen, and R. Sheroni, “Scendd: A scenario-based naturalistic driving dataset,” in *2022 IEEE 25th International Conference on Intelligent Transportation Systems (ITSC)*, pp. 4363–4368, IEEE, 2022.

- [111] 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.
- [112] 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.
- [113] J. Binas, D. Neil, S.-C. Liu, and T. Delbruck, "Ddd17: End-to-end davis driving dataset," *arXiv preprint arXiv:1711.01458*, 2017.
- [114] 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.
- [115] 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.
- [116] 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.
- [117] Z. Che, G. Li, T. Li, B. Jiang, X. Shi, X. Zhang, Y. Lu, G. Wu, Y. Liu, and J. Ye, "D²-city: a large-scale dashcam video dataset of diverse traffic scenarios," *arXiv preprint arXiv:1904.01975*, 2019.
- [118] M. Aly, "Real time detection of lane markers in urban streets," in *2008 IEEE intelligent vehicles symposium*, pp. 7–12, IEEE, 2008.
- [119] 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.
- [120] 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.
- [121] 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.
- [122] 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.
- [123] M. Pitropov, D. E. Garcia, J. Rebello, M. Smart, C. Wang, K. Czarnec, and S. Waslander, "Canadian adverse driving conditions dataset," *The International Journal of Robotics Research*, vol. 40, no. 4-5, pp. 681–690, 2021.
- [124] O. Schumann, M. Hahn, N. Scheiner, F. Weishaupt, J. F. Tilly, J. Dickmann, and C. Wöhrler, "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.
- [125] 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.
- [126] D. Temel, G. Kwon, M. Prabhushankar, and G. AlRegib, "Cure-ts: Challenging unreal and real environments for traffic sign recognition," *arXiv preprint arXiv:1712.02463*, 2017.
- [127] X. Roynard, J.-E. Deschaut, 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.
- [128] 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.
- [129] 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.
- [130] 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.
- [131] 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.
- [132] 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.
- [133] 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.
- [134] J. Zhang, X. Zou, L.-D. Kuang, J. Wang, R. S. Sherratt, and X. Yu, "Cctsd2021: a more comprehensive traffic sign detection benchmark," *Human-centric Computing and Information Sciences*, vol. 12, 2022.
- [135] Y. Pan, B. Gao, J. Mei, S. Geng, C. Li, and H. Zhao, "Semanticpos: A point cloud dataset with large quantity of dynamic instances," in *2020 IEEE Intelligent Vehicles Symposium (IV)*, pp. 687–693, IEEE, 2020.
- [136] 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.
- [137] 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.
- [138] C. Caraffi, T. Vojř, J. Trefný, 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.
- [139] 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.
- [140] 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.
- [141] 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-supervised object detection dataset for autonomous driving," *arXiv preprint arXiv:2106.11118*, 2021.
- [142] 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.
- [143] 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.
- [144] 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.
- [145] 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.
- [146] J.-L. Déziel, P. Meriaux, 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.
- [147] 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.
- [148] 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.
- [149] 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.
- [150] M. Braun, S. Krebs, F. Flohr, and D. Gavrilu, "The eurocity persons dataset: A novel benchmark for object detection. arxiv 2018," *arXiv preprint arXiv:1805.07193*.
- [151] "Tusimple: <https://github.com/tusimple/tusimple-benchmark>,"
- [152] 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.

- [153] A. Kurup and J. Bos, “Dsr: A scalable statistical filter for removing falling snow from lidar point clouds in severe winter weather,” *arXiv preprint arXiv:2109.07078*, 2021.
- [154] 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.
- [155] M. Hahner, C. Sakaridis, M. Bjelic, 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.
- [156] 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.
- [157] H. Blum, P.-E. Sarlin, J. Nieto, R. Siegwart, and C. Cadena, “The fishscapes benchmark: Measuring blind spots in semantic segmentation,” *International Journal of Computer Vision*, vol. 129, pp. 3119–3135, 2021.
- [158] 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.
- [159] 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.
- [160] A. Møgelmoose, 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.
- [161] D. Griffiths and J. Boehm, “Synthcity: A large scale synthetic point cloud,” *arXiv preprint arXiv:1907.04758*, 2019.
- [162] 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.
- [163] G. Choe, S.-H. Kim, S. Im, J.-Y. Lee, S. G. Narasimhan, and I. S. Kweon, “Ranusc: Rgb and nir urban scene dataset for deep scene parsing,” *IEEE Robotics and Automation Letters*, vol. 3, no. 3, pp. 1808–1815, 2018.
- [164] 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.
- [165] 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.
- [166] 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.
- [167] H. Schafer, E. Santana, A. Haden, and R. Biasini, “A commute in data: The comma2k19 dataset,” *arXiv preprint arXiv:1812.05752*, 2018.
- [168] 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.
- [169] 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.
- [170] Z. Wang, S. Ding, Y. Li, J. Fenn, S. Roychowdhury, A. Wallin, L. Martin, S. Rylvola, 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.
- [171] 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.
- [172] 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.
- [173] 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.
- [174] 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.
- [175] O. Jayasinghe, S. Hemachandra, D. Annettigama, 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.
- [176] 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.
- [177] 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.
- [178] 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.
- [179] “Sap: https://cs.stanford.edu/anenber/ua_data/,”
- [180] M. Meyer and G. Kuschik, “Automotive radar dataset for deep learning based 3d object detection,” in *2019 16th european radar conference (EuRAD)*, pp. 129–132, IEEE, 2019.
- [181] 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.
- [182] 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.
- [183] 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.
- [184] 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.
- [185] 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.
- [186] 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.
- [187] J. Hou, Q. Chen, Y. Cheng, G. Chen, X. Xue, T. Zeng, and J. Pu, “Suprs: 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.
- [188] 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.
- [189] 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.
- [190] L. Wang, L. Lei, H. Song, and W. Wang, “The neolix open dataset for autonomous driving,” *arXiv preprint arXiv:2011.13528*, 2020.
- [191] J. Gebele, B. Stühr, 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.
- [192] 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.
- [193] W. Zimmer, C. Creß, H. T. Nguyen, and A. C. Knoll, “A9 intersection dataset: All you need for urban 3d camera-lidar roadside perception,” *arXiv preprint arXiv:2306.09266*, 2023.

- [194] 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.
- [195] 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.
- [196] 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.
- [197] 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.
- [198] 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.
- [199] 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.
- [200] 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.
- [201] 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.
- [202] 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.
- [203] 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.
- [204] 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.
- [205] E. Arnold, S. Mozaafari, 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.
- [206] A. Kloukiniotis, A. Papandreou, C. Anagnostopoulos, A. Lalos, P. Kapsalas, D.-V. Nguyen, and K. Moustakas, "Carlasenes: 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.
- [207] 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.
- [208] 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.
- [209] 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.
- [210] 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.
- [211] 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.
- [212] A. Z. Zhu, D. Thakur, T. Özarslan, 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.
- [213] 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.
- [214] 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.
- [215] 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.
- [216] 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.
- [217] 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.
- [218] M. Wrenninge and J. Unger, "Synscapes: A photorealistic synthetic dataset for street scene parsing," *arXiv preprint arXiv:1810.08705*, 2018.
- [219] M. De Deuge, A. Quadros, C. Hung, and B. Douillard, "Unsupervised feature learning for classification of outdoor 3d scans," in *Australasian conference on robotics and automation*, vol. 2, University of New South Wales Kensington, Australia, 2013.
- [220] 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.
- [221] 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.
- [222] 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.
- [223] 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.
- [224] 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.
- [225] 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.
- [226] 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.
- [227] A. Carballo, J. Lambert, A. Monroy, 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.
- [228] 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.
- [229] 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.
- [230] 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.

- [231] 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.
- [232] 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.
- [233] 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.
- [234] 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.
- [235] 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 by distilling multimodal knowledge,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11612–11621, 2021.
- [236] 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.
- [237] 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.
- [238] 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.
- [239] 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.
- [240] T. Fleck, K. Daaboul, M. Weber, P. Schöner, 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.
- [241] K. Behrendt, “Boxy vehicle detection in large images,” in *Proceedings of the IEEE/CVF international conference on computer vision workshops*, pp. 0–0, 2019.
- [242] 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.
- [243] 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.
- [244] 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.
- [245] 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.
- [246] 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.
- [247] 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.
- [248] 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.
- [249] 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.
- [250] 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.
- [251] 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.
- [252] 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.
- [253] 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.
- [254] 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.
- [255] 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.
- [256] I. Klein, “Nexet—the largest and most diverse road dataset in the world,” *Medium*, 2017.
- [257] “Diml: <https://dimlrgb.github.io/>,”
- [258] 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.
- [259] 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.
- [260] “Flir: <https://www.flir.com/oem/adas/adas-dataset-form/>,”
- [261] K. Cordes and H. Broszio, “Vehicle lane merge visual benchmark,” in *2020 25th International Conference on Pattern Recognition (ICPR)*, pp. 715–722, IEEE, 2021.
- [262] 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.
- [263] 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.
- [264] 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.
- [265] 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.
- [266] A. Petrovai, A. D. Costea, and S. Nedeveschi, “Semi-automatic image annotation of street scenes,” in *2017 IEEE intelligent vehicles symposium (IV)*, pp. 448–455, IEEE, 2017.
- [267] 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.
- [268] 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.
- [269] 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.

- [270] 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.
- [271] 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.
- [272] 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.
- [273] 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.
- [274] 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.
- [275] N. Manikandan and K. Ganesan, "Deep learning based automatic video annotation tool for self-driving car," *arXiv preprint arXiv:1904.12618*, 2019.
- [276] D. Schörrhuber, 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.
- [277] 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.
- [278] 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.
- [279] W. Zimmer, A. Ranges, 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.
- [280] 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.
- [281] 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.
- [282] 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.
- [283] 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.
- [284] P. Krähenbühl, "Free supervision from video games," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2955–2964, 2018.
- [285] 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.
- [286] 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.
- [287] 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.
- [288] 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.
- [289] 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.
- [290] 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.
- [291] 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.
- [292] T. Deruytere, S. Vandenheide, D. Grujicic, L. Van Gool, and M.-F. Moens, "Talk2car: Taking control of your self-driving car," *arXiv preprint arXiv:1909.10838*, 2019.
- [293] 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.
- [294] 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.
- [295] 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.
- [296] D. Wu, W. Han, T. Wang, Y. Liu, X. Zhang, and J. Shen, "Language prompt for autonomous driving," *arXiv preprint arXiv:2309.04379*, 2023.
- [297] 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.
- [298] 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.
- [299] 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.
- [300] 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.
- [301] 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.
- [302] 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.
- [303] 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.
- [304] 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.
- [305] 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.
- [306] 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.
- [307] 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.
- [308] 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.
- [309] 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.
- [310] S. W. Kim, J. Phillion, 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.
- [311] 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.

- [312] 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.
- [313] 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.
- [314] 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.
- [315] Ö. Erken 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.
- [316] 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.
- [317] 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.