

Towards learning-based planning: The nuPlan benchmark for real-world autonomous driving

Napat Karnchanachari* Dimitris Geromichalos* Kok Seang Tan Nanxiang Li Christopher Eriksen
Shakiba Yaghoubi Noushin Mehdipour Gianmarco Bernasconi Whye Kit Fong Yiluan Guo Holger Caesar†

Abstract—Machine Learning (ML) has replaced handcrafted methods for perception and prediction in autonomous vehicles. Yet for the equally important planning task, the adoption of ML-based techniques is slow. We present nuPlan, the world’s first real-world autonomous driving dataset and benchmark. The benchmark is designed to test the ability of ML-based planners to handle diverse driving situations and to make safe and efficient decisions. We introduce a new large-scale dataset that consists of 1282 hours of diverse driving scenarios from 4 cities (Las Vegas, Boston, Pittsburgh, and Singapore) and includes high-quality auto-labeled object tracks and traffic light data. We mine and taxonomize common & rare driving scenarios which are used during evaluation to get fine-grained insights into the performance and characteristics of a planner. Beyond the dataset, we provide a simulation and evaluation framework that enables a planner’s actions to be simulated in closed-loop to account for interactions with other traffic participants. We present a detailed analysis of numerous baselines and investigate gaps between ML-based and traditional methods. Find the nuPlan dataset and code at nuplan.org.

I. INTRODUCTION

IN the last decade, autonomous vehicle perception and prediction have been revolutionized by deep learning-based methods trained on large-scale datasets [1], [2], [3], [4], [5], [6], [7]. While similar attempts have been made in the field of learning-based or neural planning, these are not yet able to surpass their rule-based counterparts. One possible reason is the difficulty of generalizing driving scenarios when learned from a limited number of examples. Furthermore, driving scenarios typically follow a long-tail distribution, which further exacerbates the generalization issue. Finally, learning-based planning lacks formal safety guarantees, thus making it potentially unsafe and challenging to certify.

We introduce the nuPlan dataset and simulation framework for autonomous vehicle planning. Our goal is to create a testbed for open-loop and closed-loop planning starting in real-world scenarios. This test bed is then used to compare traditional, learning-based, and hybrid planners. nuPlan enables numerous novel types of research, such as learning-based planning, the interplay between prediction and planning, and end-to-end planning using a large amount of published sensor data. We make the following contributions:

- We release the largest dataset for autonomous driving to date, with a total of 1282h from 4 cities. We also publish an unprecedented 128h of sensor data.
- We develop techniques to auto-label the dataset with accurate object tracks, traffic lights, and scenario labels.

Authors are or were (†) with Motional. (*) indicates equal contribution.

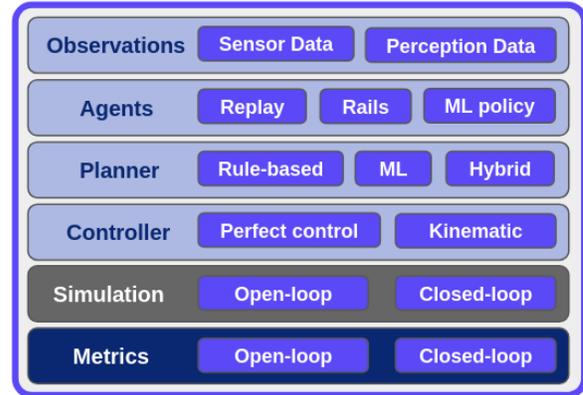


Fig. 1. An overview of the nuPlan simulation framework.

- We publish our closed-loop simulation and evaluation framework (Fig. 1) and compare the performance of traditional and learning-based planners to identify gaps.

II. RELATED WORK

In this section we discuss the most important datasets and simulators for autonomous vehicle planning. Based on these, we introduce the main categories of planners: classical, learning-based, and hybrid planners.

A. Datasets

Tab. I provides an overview of large-scale prediction and planning datasets with more than 20h of data. We omit smaller datasets like Interaction [8], highD [9], inD [10], OpenDD [11] and CommonRoad [12] and subsets of datasets not focused on prediction or planning [13], [14], [15]. With the exception of nuPlan and CommonRoad [12], all datasets focus on prediction (motion forecasting). Offline perception [16] is crucial to train and evaluate planners using high-quality object tracks, but only present in Waymo [17], MONA [18] and nuPlan. Likewise, the availability of traffic light statuses is crucial for realistic traffic simulation, but only Waymo [17] and Lyft [19] contain these and only from an online perception system, rather than developing offline traffic light status inference as in nuPlan. To assess planning performance we need to focus on specific scenarios and evaluate them in closed-loop. nuPlan is the first dataset to feature both scenario tags and a closed-loop simulation framework. Lyft [19] provides interactive tutorials for closed-loop simulation, but they lack the modular framework, evaluation server, and hold-out test set of nuPlan. Finally, we need a large-scale dataset to generalize well. Only the

Lyft [19], Shifts [15] and nuPlan datasets provide more than 1000h of driving data and nuPlan was the first such dataset that provides lidar and camera sensor data (128h), although Waymo [17] later also released compressed lidar data.

B. Simulation

Many works in the literature use proprietary simulators [20], [21]. While sharing their approaches to planning, they do not provide enough information to reproduce their experimental results. Graphical simulators like CARLA [22] and AirSim [23] focus on photorealistic rendering, but lack the realism of real-world maps and agent behavior. CommonRoad [12] was the first open-source simulator focused on planning. However, CommonRoad [12] does not provide a real world dataset as the basis for the simulation, instead resorting to a small number of manually crafted scenarios and tools to import other datasets, albeit without any sensor data. With nuPlan we aim to overcome the above limitations by releasing a large-scale real-world dataset and an open-source closed-loop simulator. Following the release of nuPlan, ScenarioNet [24] focused specifically on Reinforcement Learning and integrated nuPlan and other datasets into their pipeline. They also interfaced with MetaDrive [25] that enables graphical simulation of nuPlan.

C. Planning

Classical planning. The planning problem has long been treated as an optimization problem in the traditional approaches [26], [27], [28], [29], [30], [31]. By carefully designing a cost function, the optimization aims to generate the optimal trajectory that minimizes the cost function in the corresponding search space (e.g., A* search [27], [28], sampling-based methods [32], [33], dynamic programming [30]). While these approaches enjoy the theoretical guarantees on the convergence to an optimal solution, hand-crafting the cost function that represents the human-like driving behavior is challenging. In practice, many studies rely on tremendous engineering efforts to fine-tune the solution.

Learning based planning. Pioneering by the study of [34], the idea of using a neural network to imitate expert driver and directly output driving control command provides an alternative planning solution. With the recent success of deep learning, the learning based planning received considerable attention [35], [36], [37], [38], [39], [40], [20], [41].

Imitation learning (IL) and inverse reinforcement learning (IRL): IL trains a model to either map the sensor data directly (end-to-end system), or indirectly through the perception and prediction models (modularized system), to the expert driver actions (e.g. steering and speed profile) [34], [37], [36], [42]. With the advancements in deep learning literature, IL studies adopt the state-of-the-art supervised learning models architectures to learn better scene representations [43], [20]. IL often suffers from a poor generalization where the compounding error leads to driving scenarios that are outside of the training data, known as “covariate shift” [34]. Carefully designed data augmentation is often used to address this issue [36]. As an alternative to directly imitating the driver

behavior, IRL aims to learn an unknown reward function that explains expert demonstrations [44]. Once learned, such reward function is used to infer the optimal trajectory from a set of pre-defined or generated trajectories [45]. The maximum entropy formulation of IRL [46] has been applied to autonomous driving where the reward function is estimated based on a set of handcrafted features in [47], [48], [49].

Reinforcement learning (RL): RL learns the optimal driving behavior by interacting with the environment and optimizing a given reward function. RL is well-suited for handling the interaction between the agent and the environment in a sequential decision process. However, due to its learning by “trial-and-error” search nature, studies in RL rely on the driving simulation to provide the environment [50], [51], [52]. These studies have demonstrated strong performance in simulations. The real-world applications of RL in autonomous driving are also reported in [53], [54].

Hybrid solutions. Hybrid solutions are proposed to leverage the advantages of both the classical and learning based planning. Several studies use learning to improve the classical planning algorithm. These include using a learning model to guide the exploration for sampling-based path-planners [55], using a learning model to improve the efficiency of sampling-based motion-planners in high dimensional setting [56], applying optimizer to actively rectifies the learning model’s plan to satisfy the safety and comfort requirements [57], [40]. Other studies leverage the classical planning to generate the trajectory candidates, which are passed to an ML-based model to evaluate [49], [38], [58].

System design. While planning is the ultimate goal of the autonomous driving system, different system designs to assemble the perception and prediction introduce opportunities and challenges to improve the planning performance. Most autonomous driving systems use a multi-stage pipeline of independent tasks like perception, prediction and planning [59], [60], [61]. The hope is that the performance gain on individual tasks translates to a better planning performance. In contrast, various studies consider multi-task learning (MTL) where they jointly train models to perform perception, prediction and planning simultaneously [62], [63], [58], [38]. These works have shown that MTL achieves better data utilization at lower computation cost. Recently, some studies leverage the query-based design in transformer architectures to integrate all tasks in a unified framework that is trainable end-to-end [64], [65], [66], [67], [41], [68]. Such a framework encourages better spatio-temporal feature learning from the sensor data and directly improves the planning performance.

III. DATASET

In this section, we describe how we collect the nuPlan dataset. We enhance it by auto-labeling the object tracks of other agents and traffic lights, as well as mining for scenarios that are relevant for tracking.

A. Data collection

We collected data from 4 cities (Boston, Pittsburgh, Las Vegas, and Singapore) to build a benchmark dataset for ML-

TABLE I

AN OVERVIEW OF DATASETS FOR PREDICTION AND PLANNING (* SEE SEC. II-A FOR A DETAILED DISCUSSION)

Dataset	Year	Tasks	Perception	Traffic Lights	Scenario Tags	Closed-loop	Total Volume (h)	Sensor Volume (h)
Argoverse 1 [13]	2019	Pred	online	no	no	no	320	0
Waymo [17]	2019	Pred	offline	online	no	no	574	0*
Lyft [19]	2020	Pred	online	online	no	no*	1118	0
Shifts [15]	2021	Pred	n/a	no	no	no	1667*	0
MONA [18]	2022	Pred	offline	no	no	no	130	0
Argoverse 2 [14]	2023	Pred	online	no	no	no	763	0
nuPlan	2021	Pred+Plan	offline	offline	yes	yes	1282	128

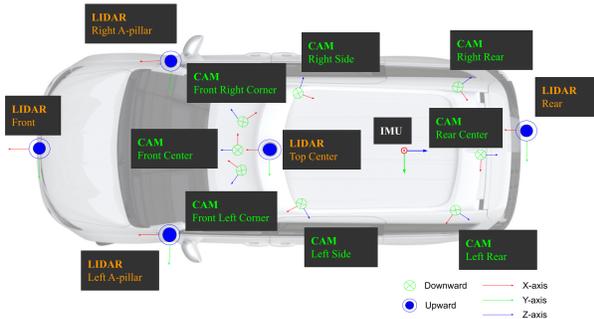


Fig. 2. The data collection vehicle’s sensor setup.

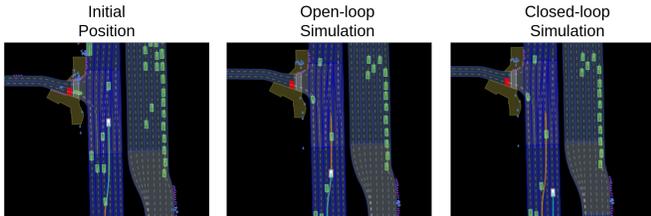


Fig. 3. Open-loop (OL) v.s. closed-loop (CL) simulation. In OL, ego vehicle follows the ground-truth trajectories. In CL, ego vehicle follows the planner model output. As shown above, starting from the initial scenario (left), using a trained machine learning planner (Urban Driver model in Sec. V), we show the different simulation roll-out in open-loop (middle) and closed-loop (right). [Ego ground-truth trajectory, Planner output trajectory]

based planning. In total, we have 1282 hours of challenging and real-world driving scenarios. For example, double parking in Boston, custom precedence patterns for left turns in Pittsburgh, crowded pick-up and drop-off points (PUDOs) in casinos in Las Vegas, and left-hand traffic in Singapore. We exclude heavy rain and night data, as these would impact the quality of our perception system (see Sec. III-B).

Manual driving. We use Chrysler Pacifica Plug-in Hybrid Electric Vehicles (PHEV) to drive in these cities. See Fig. 2 for the sensor setup. Our vehicle operators (VOs) are instructed to use a natural driving style and drive safely. Since our focus is on planning, it is crucial that we drive manually, while most other datasets [13], [14], [17], [19] use a combination of manual and automated driving, which may lead the planner to imitate less desirable driving behavior. The VOs drive from a predefined starting point to a goal using a known route. For example, we drive between various hotels and casinos on the Las Vegas strip, which are typical routes for our robotaxi and are known and mapped beforehand.

Sensor data. Sensor data include lidar point clouds and camera images. Due to the vast scale of the full sensor dataset (200+ TB), we only release a subset of the sensor data which totals 128 hours. This subset was selected to satisfy

all stratification constraints as described below.

Maps. Similar to nuScenes [69], nuPlan provides detailed human-annotated 2D high-definition semantic maps of the driving locations. We release rasterized and vectorized maps. While rasterized maps are useful for simplicity and efficient lookup, vectorized maps provide more precise geometric information and metadata. Examples of semantic map layers are lanes, car parks, crosswalks and stop lines (see Fig. 4).

B. Auto-labeling

In order to faithfully reconstruct various driving scenarios, we develop an auto-labeling system. It first generates the tracks for all the objects in the scene; then traffic light statuses are inferred from these tracks. Based on the above labels, we can reliably mine different driving scenarios.

Offline perception. We build an offline perception system to label the objects in the scene [16] automatically. Compared with the online perception systems used in many other datasets [19], [13], the offline version is not constrained by latency and causality. Therefore, the fidelity of the generated tracks is drastically higher, enabling us to evaluate planning performance under very limited perception noise.

Inspired by [16], our offline perception system contains three stages: 3D object detection, offline tracking, and global track refinement. The detector in the first stage takes the point clouds from both the top lidar and side lidars as input and detects the bounding boxes through a large neural network [70]. The offline tracker leverages both past and future detections in an extended time window to generate tracks. In the last step, a novel network is developed to load both tracks and points clouds within the tracks to refine the attributes of all the bounding boxes of vehicle class, such as positions, headings, sizes and velocity.

Traffic light status. To create a realistic simulation of the environment, it is crucial to capture the traffic light statuses. Existing datasets lack traffic lights [13], [14] or use online vision-based systems to detect their statuses [17], [19]. In contrast, we develop a novel offline system to automatically label the statuses of traffic lights by inferring them from the motion of the actors present in the scene. Our labeling system is able to cover all lanes with observed agents.

We make use of the detections and tracks produced by our offline perception system, as well as map information. To infer a green traffic light status within a given intersection, we determine if there are agents moving within the intersection in the direction controlled by the particular traffic light. To infer a red traffic light status, we check for agents slowing

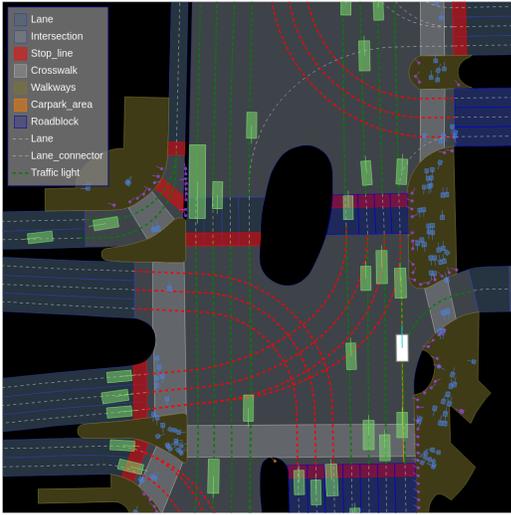


Fig. 4. Semantic map of nuPlan with 10 semantic layers in different colors, polygons and lines. It also includes traffic light statuses encoded into the lane connectors.

down or being stationary in the lane that approaches the intersection and controlled by the particular traffic light.

Scenario mining. Traditional approaches to evaluate planning performance are dominated by monotonous lane following scenarios. To get a fine-grained understanding of the planning performance, we develop a scenario taxonomy and scenario mining algorithms using low-level attributes like vehicle speed and state transitions. The attributes can be inferred from offline perception tracks and traffic light statuses. In total, we have 73 unique scenario types.

IV. SIMULATION

nuPlan provides a simulation framework (Fig. 1) that is modular and flexible to work with different datasets and setups. The simulation is initialized with the real-world *observations* captured in the dataset, namely raw sensor data or object tracks. Given these environment observations, an *agent* model can be used to predict the future trajectories of all agents. Observations and agent trajectories are passed to a *planner* that predicts the best route for the ego vehicle given the other agents' routes. Finally, a *controller* converts the intended route into a feasible trajectory. The simulation can either playback the actions recorded in the dataset (open-loop) or allow the simulation to deviate from the recording by incorporating the ego's actions (closed-loop). Below are the simulation components in detail.

A. Agents

Of the 6 object classes, vehicle, pedestrian, generic object, traffic cone, barrier, bicycle, construction zone sign, 3 are moving object classes simulated as agents: vehicles, pedestrians, and cyclists. Agents are dynamic objects that can move as the scenario evolves.

A well-known method is to simply propagate agents according to the logged data. We refer to these as non-reactive log-replay agents. Log-replay agents are used to simulate scenarios both in open-loop and closed-loop, a near-perfect recreation of the recorded data. However, closed-loop

simulation quickly diverges if the planner decides to take different actions from what is recorded in the log. Thus, in closed-loop simulation agents can be also simulated with the aim to interact with the ego vehicle and with each other by producing novel simulation states that resemble real agent behaviors. We refer to these as reactive agents. Reactive agents are only relevant in closed-loop simulation and by definition open-loop simulation uses non-reactive agents.

We develop reactive agents following the Intelligent Driver Model [71] (IDM) policy. IDM agents are initialized with the initial pose and velocity of the logged agents. The agents follow the lane center line of the underlying map. The longitudinal control is dictated by the IDM policy. This allows the agents to react to the ego's actions, as well as other reactive agents. In turn, this reduces false collisions and lets scenarios play out for longer. Note that we only apply this policy to vehicles, while Vulnerable Road Users (VRUs) are replayed from the log. We choose not to model VRUs as reactive agents as their behavior is often uncooperative and thus hard to model.

B. Controller

In nuPlan, planners provide a trajectory as a sequence of poses in $SE(2)$, without any kinematic feasibility requirement. This trajectory is assumed to be sampled at specific times in the future according to the simulation configuration. To assert kinematic feasibility and prevent users from cheating, we require the use of a controller. nuPlan provides the flexibility to use any controller, such as perfect tracking, which simply interpolates poses along the planned trajectory.

We developed a two-stage controller to propagate the simulation in closed-loop. This controller consists of two parts, a trajectory tracker and a motion model to forward-integrate the simulation. We implement a Linear Quadratic Regulator (LQR) [72], [73] as the tracker. The found optimal control policy is then fed to the second part of the simulation controller, a kinematic bicycle model which is forward-integrated to propagate the simulation state. Alternatively, we also support different trackers in the two-stage controller, such as an iterative-LQR tracker.

C. Evaluation

Different metrics and frameworks have been explored for scoring models in prediction and motion planning benchmarks [17], [74], [75], [76], [77], [78], [79], [80]. In this paper, we select a set of metrics and design an aggregation method to compare the performance of planners. In open-loop, we only evaluate the closeness of the planner generated trajectory to the human-driven trajectory. The open-loop metrics are modified Average Displacement Error (ADE), Final Displacement Error (FDE), Average Heading Error (AHE), Final Heading Error (FHE), and Miss Rate (MR) in which we calculate the metrics over different horizons and report their average score. For closed-loop, we use a combination of metrics to evaluate lawfulness and compliance with traffic rules consisting of no at-fault collisions, trajectories inside drivable area, no trajectories in lanes belonging to oncoming

TABLE II
PLANNER METRICS THEIR AND WEIGHTS

Simulation	Metric name	Multiplier weight	Average weight
Open-loop	MR within bound	{0, 1}	-
	AHE and FHE within bound	-	2
	ADE and FDE within bound	-	1
Closed-loop	No at-fault collisions	{0, 0.5, 1}	-
	Drivable area compliance	{0, 1}	-
	Making progress	{0, 1}	-
	Driving direction compliance	{0, 0.5, 1}	-
	TTC within bound	-	5
	Progress along route ratio	-	5
	Speed limit compliance	-	4
	Comfort	-	2

traffic, not driving above the speed limit and maintaining enough Time To Collision (TTC) with other road users, metrics to evaluate progress towards the goal and measure the rider comfort. All metrics’ scores are normalized to the range $[0 - 1]$ using thresholds that are selected based on legal requirements and natural human driving. A higher score indicates a better performance.

The final score of a planner is computed by averaging the scores for its generated trajectories across all scenarios. The score of a trajectory in a scenario is given by a hybrid weighted average of all metrics’ scores.

The rest of the metrics are weighted according to their importance (See Tab. II) and then averaged to compute the scenario score as:

$$\text{scenario score} = \prod_{i \in \text{multiplier metrics}} \text{score}_i \times \sum_{j \in \text{average metrics}} \text{weight}_j \times \text{score}_j$$

We define the score for each challenge (open-loop, closed-loop non-reactive, and closed-loop reactive) as the average scenario score across all scenarios for that challenge.

V. EXPERIMENTS

Here we present a number of planning baselines and their results when evaluated on the nuPlan benchmark. We analyze how the planning performance is impacted by lower quality perception inputs, as well as how it generalizes to other cities. Finally, we discuss the new state-of-the-art set by the submissions to the first nuPlan challenge.

A. Planning baselines

We implement several planning methods that are representative of the literature.

a) Simple Planner: The Simple planner has little planning capability. The planner plans a straight line at a constant speed. The only logic of this planner is to decelerate if the current velocity exceeds the max velocity.

b) IDM Planner: The Intelligent Driver Model (IDM) planner is essentially an Adaptive Cruise Control (ACC) policy [81]. The planner consists of two parts: path planning and longitudinal control. The path planning component is a breadth-first search algorithm. It finds a center-line path toward the mission goal extracted from the underlying map structure. The longitudinal control follows the IDM policy. The policy describes how fast the planner should go based on the distance between itself and the closest leading agent.

TABLE III
MAIN RESULTS

Planner	Open-loop	Closed-loop Non-reactive	Closed-loop Reactive
Simple Planner	0.22	0.32	0.37
IDM Planner	0.30	0.73	0.76
Raster Planner	0.52	0.47	0.46
UrbanDriver	0.90	0.68	0.67

c) Raster ML planner: Similar to the encoder in [36], [82], the raster planner uses ResNet-50 [83] as the backbone to encode features from an ego-centric multi-channel raster representing the ego, the agents and the map. The model directly outputs the final ego trajectory. The planner does not perform any post-processing on the predicted ego trajectory.

d) UrbanDriver ML Planner: We adopted an open-loop training variant of the UrbanDriver model [84] as a representative machine learning planner baseline. The model processes vectorized agents and map inputs into local feature descriptors that are passed to a global attention mechanism for yielding a predicted ego trajectory. We train the model using imitation learning to match expert trajectories available in the nuPlan dataset. Data augmentation is additionally performed on the agents and expert trajectory provided during training to mitigate data distribution drift encountered during closed-loop simulation. This version was used for the challenge. We also implemented a multi-step prediction baseline variant as discussed in [84] and originally proposed in [85] to further address the distribution shift, for the experiments in this work but do not open-source this implementation.

B. Main results

Tab. III shows the planning results for the proposed baselines in each of the three challenge setups. Supervised learning-based planners excel in an open-loop setting. This is unsurprising as the task is akin to the traditional motion forecasting challenge. This suggests that an ML planner can choose to make similar decisions to a human driver in open-loop settings. However, ML planners still struggle to overcome the distribution shift in closed-loop. A closed-loop scenario can develop into a new situation that was never present in the training dataset. Even techniques such as data augmentation and closed-loop training fail to overcome this domain gap. This is evident in both the literature [36] and our experiments. Rule-based planners, on the other hand, face no such issues. Policies like IDM can produce decent driving behavior. This is confirmed by the metrics as it achieved the highest scores for closed-loop. It should be noted though that the reactive agents are also modelled with a similar IDM. The use of similar assumptions on the vehicle behavior may result in giving the IDM planner an unfair advantage over other planners in closed-loop evaluation. It is evident that sufficiently sophisticated rule-based planners still outperform purely learned planners in closed-loop settings.

C. Perturbation

As the nuPlan dataset is created with offline perception, to capture the original probability distribution of data collected online, we injected uniform noise on the detections. Noise

TABLE IV

DETECTION PERTURBATION IN CLOSED-LOOP REACTIVE SIMULATION

Simulation Agents	UrbanDriver	UrbanDriverOnline
Original	0.67	0.60
Perturbed	0.64	0.59

TABLE V

URBANDRIVER LOCATION GENERALIZATION

Locations	Open-Loop	Closed-Loop Non-reactive	Closed-Loop Reactive
Las Vegas	0.91	0.57	0.57
Singapore	0.28	0.23	0.18
Boston	0.57	0.49	0.50
Pittsburgh	0.41	0.39	0.32

was added in the dimensions and the pose of the detected agents, with variance extracted by comparing offline and online detections. The scores of planners under nominal and noise-injected simulations in closed-loop reactive mode are presented in Tab. IV. A version of UrbanDriver trained on the perturbed data is called UrbanDriverOnline, which shows a performance deterioration compared to the nominal model on both nominal and injected data. This indicates the value of high-quality offline annotations in the dataset and the learning pipeline.

D. Generalization

The location generalization experiment is shown in Tab. V. The experiment aims to test a model’s generalization capabilities. The UrbanDriver model was trained purely on data from Las Vegas. The model was tested separately on scenarios from Singapore, Boston, Pittsburgh, and Las Vegas. The open-loop performance dropped by 53.8%, while closed-loop non-reactive and closed-loop reactive performance dropped by 35.1% and 41.5% respectively. The worst-performing location is Singapore. This can be explained by the left-hand traffic, while the model was trained on right-hand traffic. One insight is that the correlation between the model’s open-loop and closed-loop performance is relatively weak. The difference between open and closed-loop scores across Singapore, Boston, and Pittsburgh is only 16.3%, while for Las Vegas it is more than double at 37.3%. This indicates that a good motion forecasting model does not translate to closed-loop capabilities. Thus a major challenge is to overcome the domain gap between open-loop and closed-loop before tackling larger generalization problems.

E. nuPlan challenge

In the nuPlan motion planning challenge contestants create a planner to traverse a set of diverse and challenging scenarios across all four cities. Tab. VI shows the Overall Score, which is the average across Open-loop, Closed-loop Non-reactive, and Closed-loop Reactive challenges of the top four planners. In the open-loop challenge, planners that incorporated supervised learned methods scored relatively well. In the closed-loop challenges, planners employed a combination of learned and handcrafted components. A common theme was the use of a learned model to first predict the ego’s planned trajectory. [82] uses a raster-based model

TABLE VI

NUPLAN CHALLENGE LEADERBOARD

Team name	Approach	Overall Score
CS Tu [87]	Rule-based + ML refinement	0.895
AutoHorizon [82]	ML + optimizer	0.875
Pegasus	ML + collision checking	0.848
AID [86]	ML + hierarchical game theory	0.829

that outputs a spatial-temporal heatmap for the ego and an occupancy map for the surrounding agents. [86] are vector-based using transformers as a backbone. Once the trajectory is obtained, the planner has a further refinement stage used to ensure kinematic feasibility and collision avoidance. The highest-scoring planner in closed-loop was mostly rule-based [87]. It generates a handful of trajectories by perturbing the center line laterally at different velocities. Trajectories are selected with a heuristic that considers factors such as collision, drivable areas, traffic laws, and comfort. An ML-generated trajectory is fused to correct the long-term planned horizon. This limited the influence of the learned model. We draw two conclusions from the challenge results. First, ML-based methods require additional post-processing for closed-loop driving. Second, hybrid methods appear to be the most effective approach, combining traditional and data-driven methods.

VI. CONCLUSION

We presented nuPlan, the first real-world driving benchmark and the largest existing labeled autonomous driving dataset. The dataset consists of 1282 hours of diverse driving scenarios across 4 cities as well as an unprecedented 128 hours of raw sensor data and is accompanied by an evaluation framework powered by a closed-loop simulator; the dataset and the evaluation framework are publicly available. We investigated the state of current rule-based and learned-based planners by evaluating multiple approaches on the nuPlan dataset across challenging driving scenarios. The first public nuPlan challenge demonstrated that rule-based planners outperform purely ML-based ones, but hybrid planners with learned-based components show the most promise in handling difficult scenarios. In the future, we plan to mine for richer long-tail driving scenarios, design scenario-based metrics, provide ML-based planning and agent baselines and explore end-to-end planner training directly from sensor data.

ACKNOWLEDGEMENT

We would like to thank the numerous current and former colleagues at Motional who contributed to this work, especially Juraj Kabzan, Edouard Capellier, Abhinav Rai, Xiaoli Meng, Jiong Yang, Lubing Zhou, Abirami Srinivasan, Bing Jui Ho, Hiok Hian Ong, Mitchell Spryn, Taufik Tirtosudiro, Michael Noronha, Vijay Govindarajan, Jeff Wang, Nelly Lyu, Samson Hang, Shashank Chaudhary, Patrick Weygand, Vern Jensen, Abhimanyu Singh, Qiang Xu and Sammy Omari. Furthermore, many external advisors gave useful feedback, including Kashyap Chitta and Oscar de Groot.

REFERENCES

- [1] Y. Zhou and O. Tuzel, "Voxelnet: End-to-end learning for point cloud based 3d object detection," in *CVPR*, 2018.
- [2] A. H. Lang, S. Vora, H. Caesar, L. Zhou, J. Yang, and O. Beijbom, "Pointpillars: Fast encoders for object detection from point clouds," in *CVPR*, 2019.
- [3] S. Vora, A. H. Lang, B. Helou, and O. Beijbom, "Pointpainting: Sequential fusion for 3d object detection," in *CVPR*, 2020.
- [4] Q. Chen, S. Vora, and O. Beijbom, "Polarstream: Streaming object detection and segmentation with polar pillars," in *NeurIPS*, 2021.
- [5] T. Yin, X. Zhou, and P. Krähnbühl, "Center-based 3d object detection and tracking," in *CVPR*, 2021.
- [6] T. Phan-Minh, E. C. Grigore, F. A. Boulton, O. Beijbom, and E. M. Wolff, "Covernet: Multimodal behavior prediction using trajectory sets," in *CVPR*, 2020.
- [7] M. Liang, B. Yang, R. Hu, Y. Chen, R. Liao, S. Feng, and R. Urtasun, "Learning lane graph representations for motion forecasting," in *ECCV*, 2020.
- [8] W. Zhan, L. Sun, D. Wang, H. Shi, A. Clause, M. Naumann, J. Kümmerle, H. Königshof, C. Stiller, A. de La Fortelle, and M. Tomizuka, "INTERACTION dataset: An international, adversarial and cooperative motion dataset in interactive driving scenarios with semantic maps," *arXiv preprint arXiv:1910.03088*, 2019.
- [9] R. Krajewski, J. Bock, L. Kloeker, and L. Eckstein, "The highd dataset: A drone dataset of naturalistic vehicle trajectories on german highways for validation of highly automated driving systems," in *ITSC*, 2018.
- [10] 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 *IV*, 2020.
- [11] A. Breuer, J.-A. Termöhlen, S. Homoceanu, and T. Fingscheidt, "openDD: A large-scale roundabout drone dataset," in *ITSC*, 2020.
- [12] M. Althoff, M. Koschi, and S. Manzingler, "Commonroad: Composable benchmarks for motion planning on roads," in *IV*, 2017.
- [13] M.-F. Chang, J. W. Lambert, P. Sangkloy, J. Singh, S. Bak, A. Hartnett, D. Wang, P. Carr, S. Lucey, D. Ramanan, and J. Hays, "Argoverse: 3d tracking and forecasting with rich maps," in *CVPR*, 2019.
- [14] 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," in *NeurIPS Track on Datasets and Benchmarks*, 2021.
- [15] A. Malinin, N. Band, Y. Gal, M. Gales, A. Ganshin, G. Chesnokov, A. Noskov, A. Ploskonosov, L. Prokhorenkova, I. Provlkov, V. Raina, V. Raina, D. Roginskiy, M. Shmatova, P. Tigas, and B. Yangel, "Shifts: A dataset of real distributional shift across multiple large-scale tasks," in *NeurIPS Track on Datasets and Benchmarks*, 2021.
- [16] C. R. Qi, Y. Zhou, M. Najibi, P. Sun, K. Vo, B. Deng, and D. Anguelov, "Offboard 3d object detection from point cloud sequences," in *CVPR*, 2021.
- [17] S. Ettinger, S. Cheng, and B. C. et al., "Large scale interactive motion forecasting for autonomous driving: The Waymo Open Motion Dataset," in *ICCV*, 2021.
- [18] L. Gressenbuch, K. Esterle, T. Kessler, and M. Althoff, "Mona: The munich motion dataset of natural driving," in *ITSC*, 2022.
- [19] J. Houston, G. Zuidhof, and L. B. et al., "One thousand and one hours: Self-driving motion prediction dataset," in *CoRL*, 2020.
- [20] O. Scheel, L. Bergamini, M. Wolczyk, B. Osiński, and P. Ondruska, "Urban driver: Learning to drive from real-world demonstrations using policy gradients," in *CoRL*, 2022.
- [21] M. Bansal, A. Krizhevsky, and A. Ogale, "Chauffeurnet: Learning to drive by imitating the best and synthesizing the worst," in *RSS*, 2019.
- [22] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun, "CARLA: An open urban driving simulator," *CoRR*, 2017.
- [23] S. Shah, D. Dey, C. Lovett, and A. Kapoor, "Airsim: High-fidelity visual and physical simulation for autonomous vehicles," in *Field and Service Robotics: Results of the 11th International Conference*, 2018.
- [24] Z. A. Daniels and D. Metaxas, "Scenarionet: An interpretable data-driven model for scene understanding," in *IJCAI Workshop on Explainable Artificial Intelligence (XAI) 2018*, 2018.
- [25] Q. Li, Z. Peng, L. Feng, Q. Zhang, Z. Xue, and B. Zhou, "Metadrive: Composing diverse driving scenarios for generalizable reinforcement learning," *PAMI*, 2022.
- [26] B. Paden, M. Cáp, 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*, 2016.
- [27] M. Buehler, K. Iagnemma, and S. Singh, *The DARPA urban challenge: autonomous vehicles in city traffic*. Springer, 2009, vol. 56.
- [28] Z. Ajanovic, B. Lacevic, B. Shyrokau, M. Stolz, and M. Horn, "Search-based optimal motion planning for automated driving," in *IROS*, 2018.
- [29] T. Fraichard and C. Laugier, "Path-velocity decomposition revisited and applied to dynamic trajectory planning," in *ICRA*, 1993.
- [30] H. Fan, F. Zhu, C. Liu, L. Zhang, L. Zhuang, D. Li, W. Zhu, J. Hu, H. Li, and Q. Kong, "Baidu apollo EM motion planner," *arXiv preprint arXiv:1807.08048*, 2018.
- [31] M. Montemerlo, J. Becker, S. Bhat, H. Dahlkamp, D. Dolgov, S. Ettinger, D. Haehnel, T. Hilden, G. Hoffmann, B. Huhnke *et al.*, "Junior: The stanford entry in the urban challenge," *Journal of Field Robotics*, 2008.
- [32] J. J. Kuffner and S. M. LaValle, "Rrt-connect: An efficient approach to single-query path planning," in *ICRA*, 2000.
- [33] S. Karaman and E. Frazzoli, "Sampling-based algorithms for optimal motion planning," *IJRR*, 2011.
- [34] D. A. Pomerleau, "Alvinn: An autonomous land vehicle in a neural network," in *NeurIPS*, 1988.
- [35] F. Codevilla, E. Santana, A. M. López, and A. Gaidon, "Exploring the limitations of behavior cloning for autonomous driving," in *ICCV*, 2019.
- [36] M. Bansal, A. Krizhevsky, and A. Ogale, "Chauffeurnet: Learning to drive by imitating the best and synthesizing the worst," in *RSS*, 2019.
- [37] M. Bojarski, D. Del Testa, D. Dworakowski, B. Firner, B. Flepp, P. Goyal, L. D. Jackel, M. Monfort, U. Muller, J. Zhang *et al.*, "End to end learning for self-driving cars," *arXiv preprint arXiv:1604.07316*, 2016.
- [38] W. Zeng, W. Luo, S. Suo, A. Sadat, B. Yang, S. Casas, and R. Urtasun, "End-to-end interpretable neural motion planner," in *CVPR*, 2019.
- [39] J. Hawke, R. Shen, C. Gurau, S. Sharma, D. Reda, N. Nikolov, P. Mazur, S. Micklethwaite, N. Griffiths, A. Shah *et al.*, "Urban driving with conditional imitation learning," in *ICRA*, 2020.
- [40] M. Vitelli, Y. Chang, Y. Ye, A. Ferreira, M. Wolczyk, B. Osiński, M. Niendorf, H. Grimmert, Q. Huang, A. Jain *et al.*, "SafetyNet: Safe planning for real-world self-driving vehicles using machine-learned policies," in *ICRA*, 2022.
- [41] 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 *CVPR*, 2023.
- [42] M. Hallgarten, M. Stoll, and A. Zell, "From prediction to planning with goal conditioned lane graph traversals," *arXiv preprint arXiv:2302.07753*, 2023.
- [43] M. Müller, A. Dosovitskiy, B. Ghanem, and V. Koltun, "Driving policy transfer via modularity and abstraction," in *CoRL*, 2018.
- [44] P. Abbeel and A. Y. Ng, "Apprenticeship learning via inverse reinforcement learning," in *ICML*, 2004.
- [45] J. Ho and S. Ermon, "Generative adversarial imitation learning," *NIPS*, 2016.
- [46] N. Aghasadeghi and T. Bretl, "Maximum entropy inverse reinforcement learning in continuous state spaces with path integrals," in *IROS*, 2011.
- [47] M. Wulfmeier, P. Ondruska, and I. Posner, "Maximum entropy deep inverse reinforcement learning," *arXiv preprint arXiv:1507.04888*, 2015.
- [48] B. D. Ziebart, A. L. Maas, J. A. Bagnell, A. K. Dey *et al.*, "Maximum entropy inverse reinforcement learning," in *AAAI*, 2008.
- [49] T. Phan-Minh, F. Howington, T.-S. Chu, S. U. Lee, M. S. Tomov, N. Li, C. Dicle, S. Fidler, F. Suarez-Ruiz, R. Beaudoin *et al.*, "DriveIRL: Drive in real life with inverse reinforcement learning," in *ICRA*, 2023.
- [50] S. Shalev-Shwartz, S. Shammah, and A. Shashua, "Safe, multi-agent, reinforcement learning for autonomous driving," *arXiv preprint arXiv:1610.03295*, 2016.
- [51] D. Chen, B. Zhou, V. Koltun, and P. Krähnbühl, "Learning by cheating," in *CoRL*, 2020.
- [52] D. Chen, V. Koltun, and P. Krähnbühl, "Learning to drive from a world on rails," in *ICCV*, 2021.
- [53] M. Riedmiller, M. Montemerlo, and H. Dahlkamp, "Learning to drive a real car in 20 minutes," in *2007 Frontiers in the Convergence of Bioscience and Information Technologies*, 2007.
- [54] A. Kendall, J. Hawke, D. Janz, P. Mazur, D. Reda, J.-M. Allen, V.-D. Lam, A. Bewley, and A. Shah, "Learning to drive in a day," in *ICRA*, 2019.

- [55] O. Arslan and P. Tsiotras, "Machine learning guided exploration for sampling-based motion planning algorithms," in *IROS*, 2015.
- [56] A. H. Qureshi, A. Simeonov, M. J. Bency, and M. C. Yip, "Motion planning networks," in *ICRA*, 2019.
- [57] H. Pulver, F. Eiras, L. Carozza, M. Hawasly, S. V. Albrecht, and S. Ramamoorthy, "Pilot: Efficient planning by imitation learning and optimisation for safe autonomous driving," in *IROS*, 2021.
- [58] S. Casas, A. Sadat, and R. Urtasun, "Mp3: A unified model to map, perceive, predict and plan," in *CVPR*, 2021.
- [59] L. Chen, L. Platinsky, S. Speichert, B. Osiński, O. Scheel, Y. Ye, H. Grimmer, L. Del Pero, and P. Ondruska, "What data do we need for training an av motion planner?" in *ICRA*, 2021.
- [60] J. Gu, C. Sun, and H. Zhao, "Densent: End-to-end trajectory prediction from dense goal sets," in *ICCV*, 2021.
- [61] K. Xiong, S. Gong, X. Ye, X. Tan, J. Wan, E. Ding, J. Wang, and X. Bai, "Cape: Camera view position embedding for multi-view 3d object detection," in *CVPR*, 2023.
- [62] A. Sadat, S. Casas, M. Ren, X. Wu, P. Dhawan, and R. Urtasun, "Perceive, predict, and plan: Safe motion planning through interpretable semantic representations," in *ECCV*, 2020.
- [63] J. Ngiam, V. Vasudevan, B. Caine, Z. Zhang, H.-T. L. Chiang, J. Ling, R. Roelofs, A. Bewley, C. Liu, A. Venugopal *et al.*, "Scene transformer: A unified architecture for predicting future trajectories of multiple agents," in *ICLR*, 2022.
- [64] D. Chen and P. Krähenbühl, "Learning from all vehicles," in *CVPR*, 2022.
- [65] S. Casas, A. Sadat, and R. Urtasun, "MP3: A unified model to map, perceive, predict and plan," in *CVPR*, 2021.
- [66] K. Chitta, A. Prakash, and A. Geiger, "Neat: Neural attention fields for end-to-end autonomous driving," in *ICCV*, 2021.
- [67] S. Hu, L. Chen, P. Wu, H. Li, J. Yan, and D. Tao, "St-p3: End-to-end vision-based autonomous driving via spatial-temporal feature learning," in *ECCV*, 2022.
- [68] P. Wu, X. Jia, L. Chen, J. Yan, H. Li, and Y. Qiao, "Trajectory-guided control prediction for end-to-end autonomous driving: A simple yet strong baseline," in *NeurIPS*, 2022.
- [69] 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 *CVPR*, 2020.
- [70] I. Radosavovic, R. P. Kosaraju, R. Girshick, K. He, and P. Dollár, "Designing network design spaces," in *CVPR*, 2020.
- [71] M. Treiber, A. Hennecke, and D. Helbing, "Congested traffic states in empirical observations and microscopic simulations," *Phys. Rev. E*, 2000.
- [72] J. Liu, Z. Yang, Z. Huang, W. Li, S. Dang, and H. Li, "Simulation performance evaluation of pure pursuit, stanley, lqr, mpc controller for autonomous vehicles," in *2021 IEEE international conference on real-time computing and robotics (RCAR)*, 2021.
- [73] B. Varma, N. Swamy, and S. Mukherjee, "Trajectory tracking of autonomous vehicles using different control techniques (pid vs lqr vs mpc)," in *2020 International conference on smart technologies in computing, electrical and electronics (ICSTCEE)*, 2020.
- [74] M. Hekmatnejad, S. Yaghoubi, A. Dokhanchi, H. B. Amor, A. Shrivastava, L. Karam, and G. Fainekos, "Encoding and monitoring responsibility sensitive safety rules for automated vehicles in signal temporal logic," in *International Conference on Formal Methods and Models for System Design*, 2019.
- [75] W. Xiao, N. Mehdipour, A. Collin, A. Y. Bin-Nun, E. Frazzoli, R. D. Tebbens, and C. Belta, "Rule-based optimal control for autonomous driving," in *Proceedings of the ACM/IEEE 12th International Conference on Cyber-Physical Systems*, 2021, pp. 143–154.
- [76] J. Houston, G. Zuidhof, L. Bergamini, Y. Ye, L. Chen, A. Jain, S. Omari, V. Igllovikov, and P. Ondruska, "One thousand and one hours: Self-driving motion prediction dataset," in *CoRL*, 2021.
- [77] R. McAllister, B. Wulfe, J. Mercat, L. Ellis, S. Levine, and A. Gaidon, "Control-aware prediction objectives for autonomous driving," in *ICRA*, 2022.
- [78] S. Maierhofer, P. Moosbrugger, and M. Althoff, "Formalization of intersection traffic rules in temporal logic," in *IV*, 2022.
- [79] S. Jiang, Z. Xiong, W. Lin, Y. Cao, Z. Xia, J. Miao, and Q. Luo, "An efficient framework for reliable and personalized motion planner in autonomous driving," *RA-L*, 2022.
- [80] M. Ilievski, "Wisebench: A motion planning benchmarking framework for autonomous vehicles," Master's thesis, University of Waterloo, 2020.
- [81] O. Derbel, T. Peter, H. Zebiri, B. Mourllion, and M. Basset, "Modified intelligent driver model for driver safety and traffic stability improvement," *IFAC Proceedings Volumes*, vol. 46, no. 21, pp. 744–749, 2013.
- [82] Y. Hu, K. Li, P. Liang, J. Qian, Z. Yang, H. Zhang, W. Shao, Z. Ding, W. Xu, and Q. Liu, "Imitation with spatial-temporal heatmap: 2nd place solution for nuplan challenge," 2023.
- [83] R. Wightman, "Pytorch image models," <https://github.com/rwightman/pytorch-image-models>, 2019.
- [84] O. Scheel, L. Bergamini, M. Wolczyk, B. Osinski, and P. Ondruska, "Urban driver: Learning to drive from real-world demonstrations using policy gradients," in *CoRL*, 2021.
- [85] A. Venkatraman, M. Hebert, and J. Bagnell, "Improving multi-step prediction of learned time series models," *AAAI*, 2015.
- [86] Z. Huang, H. Liu, and C. Lv, "Gameformer: Game-theoretic modeling and learning of transformer-based interactive prediction and planning for autonomous driving," 2023.
- [87] D. Dauner, M. Hallgarten, A. Geiger, and K. Chitta, "Parting with misconceptions about learning-based vehicle motion planning," 2023.
- [88] Y. Zhou, P. Sun, Y. Zhang, D. Anguelov, J. Gao, T. Ouyang, J. Guo, J. Ngiam, and V. Vasudevan, "End-to-end multi-view fusion for 3d object detection in lidar point clouds," in *CoRL*, 2020.
- [89] Y. Wang, A. Fathi, A. Kundu, D. A. Ross, C. Pantofaru, T. Funkhouser, and J. Solomon, "Pillar-based object detection for autonomous driving," in *ECCV*, 2020.
- [90] X. Weng, J. Wang, D. Held, and K. Kitani, "3d multi-object tracking: A baseline and new evaluation metrics," in *IROS*, 2020.
- [91] B. Yang, M. Bai, M. Liang, W. Zeng, and R. Urtasun, "Auto4d: Learning to label 4d objects from sequential point clouds," *arXiv preprint arXiv:2101.06586*, 2021.
- [92] J. Redmon and A. Farhadi, "Yolov3: An incremental improvement," *arXiv preprint arXiv:1804.02767*, 2018.
- [93] S. R. Rath, "Traffic light detection using yolov3," <https://github.com/sovit-123/Traffic-Light-Detection-Using-YOLOv3/>, 2020.
- [94] M. B. Jensen and M. Philip, "Lisa traffic light dataset," <https://www.kaggle.com/datasets/mbornoe/lisa-traffic-light-dataset/>, 2015.
- [95] M. Igl, D. Kim, A. Kuefler, P. Mougín, P. Shah, K. Shiarlis, D. Anguelov, M. Palatucci, B. White, and S. Whiteson, "Symphony: Learning realistic and diverse agents for autonomous driving simulation," in *ICRA*, 2022.

SUPPLEMENTARY MATERIAL

In this supplementary material we provide additional information about the creation of the nuPlan dataset, the simulation framework and evaluation protocol, as well as additional experiments.

VII. DATASET

A. Data collection

Here we describe the sensor setup, synchronization, dataset splits, semantic map layers and object annotation statistics.

Sensor setup. The nuPlan sensor setup is collected by a fleet of 31 vehicles with identical sensor setup (Tab. VII). The vehicles are equipped with 5 lidars, 8 cameras, as well as GNSS & IMU. Thus both lidars and cameras cover the full 360-degree environment and minimize blindspots (Fig. 5). The lidars can generate lidar point clouds with up to 360k and 720k points for 20-channel and 40-channel lidars, respectively.

Synchronization. To ensure high-quality data alignment between multiple sensors, lidar sensors are synchronized with the system time on the car in multiple lidar periods. The exposure of each camera is triggered in the same way plus a camera-specific offset, where the offset is an optimal value to ensure images are acquired right when the lidars are sweeping through each camera’s FOV. Given there are multiple lidar sweeps in a spin, we merge a full sweep from each lidar sensor into a merged sweep since they complete a revolution at nearly the same instance. After they are merged, the merged point clouds are transformed to the same frame coordinate and timestamp. We also perform ego motion compensation to take ego position into account when the points are acquired.

Dataset splits. The dataset splits (train, val, and test) are geographically overlapping, which means they can cover the same region in a city. However, to minimize temporal data leakage, data from the same day *and* city is not shared across splits. In addition, the dataset is stratified across days, cities, and driving scenarios to ensure all splits have similar balances across these dimensions.



Fig. 5. Examples of nuPlan sensor data. We can see the lidar data and map in the center and 6 camera views around it. All views show the highly accurate autolabeling annotations.

TABLE VII
SENSOR SETUP IN NUPLAN

Sensor	Make	Freq.	Details
1x Top Lidar	Pandar 40P	20Hz	Spinning, 40 channels, 360° horizontal FOV, $\leq 200\text{m}$ range, $\leq \pm 2\text{cm}$, 40° ([-25°, 15°]) vertical FOV, $\leq 720\text{k}$ points per second
2x A-Pillar Lidars	Pandar 40P	20Hz	Spinning, 40 channels, 360° horizontal FOV, $\leq 200\text{m}$ range, $\leq \pm 2\text{cm}$, 40° ([-25°, 15°]) vertical FOV, $\leq 720\text{k}$ points per second
2x Bumper Lidars	Pandar 20PA front, 20PB rear	20Hz	Spinning, 20 channels, 360° horizontal FOV, $\leq 200\text{m}$ range, $\leq \pm 2\text{cm}$, 33° ([-25°, 8°]) front vertical FOV, 22° ([-19°, 3°]) rear vertical FOV, $\leq 360\text{k}$ points per second
8x Cameras	D3 Engineering D3RCM	10Hz	RGB, 1/2.7" CMOS sensor, 2000x1200 resolution, split-pixel image sensor
1x GNSS	Trimble BX992	20Hz	Position latency $< 20\text{ms}$, 0.10° roll/pitch
1x IMU	Honeywell HG1120	100Hz	MEMS gyroscopes, accelerometers and magnetometers

Map layers. We define the semantic map layers in nuPlan. The map layers are designed with planning in mind. A graph can be constructed from the network of lane and lane connectors. Additionally, the map is annotated with important road features (stop lines, crosswalks, car parks) that impact the decision-making of the AV. See Tab. VIII for more details.

Annotation statistics. Here we present more statistics for the annotations in nuPlan. Fig. 6 shows the number of tracks for each of the 6 object classes in nuPlan. While there are about 2 orders of magnitude between the most (generic object) and least (construction zone sign) common class, the dataset is nevertheless less long-tailed than other datasets with more fine-grained classes [69]. Furthermore, even the rarest class still has more than 10^5 tracks, which shows the advantage of such a large dataset.

We also present size statistics for different classes in Fig. 7. These show that our boxes are statistically stable in their sizes except for the vehicle class that includes construction vehicles with irregular shapes.

Fig. 8 shows the absolute velocities for three classes (vehicle, pedestrian, and bicycle), since other classes are static most of the time. We observe that most of the objects are slow-moving, which represents primarily urban scenarios, while a small number of pedestrians and bicycles are moving at unrealistically high speeds, possibly due to noisy annotations.

Fig. 9 shows the object (box) orientations relative to the ego vehicle orientation. Due to the grid-like road layout, most vehicles have orientations that are multiples of $\pm 90^\circ$.

TABLE VIII
SEMANTIC MAP LAYERS IN NUPLAN

Layer name	Layer description
Baseline paths	The center line along lanes and lane connectors.
Carpark areas	Polygons representing car parks.
Crosswalks	Polygons representing crosswalks.
Generic drivable areas	Polygons representing areas where the AV is allowed to drive.
Intersections	An area connecting multiple road segments together.
Lane connectors	Road segments connecting two lanes. Typically, they exist within an intersection.
Lane group connectors	A group of adjacent lane connectors that travel in the same direction
Lane groups polygons	A group of adjacent lanes that travel in the same direction.
Lanes polygons	Polygons representing a lane.
Road segments	Groupings of lane groups that travel in opposite directions from one another.
Stop polygons	Polygon in which the AV may be required to stop.
Traffic lights	3D locations of a set of traffic lights.
Walkways	Polygons representing walkways.

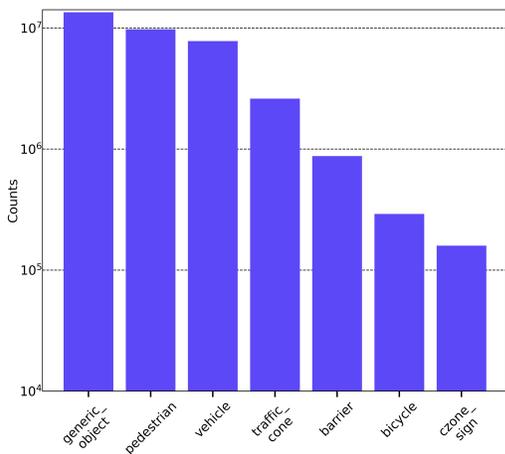


Fig. 6. Number of tracks per category.

Fig. 10 shows the spatial coverage of our ego vehicles across all maps. In Las Vegas, most routes start and end in PUDOs, which leads to these areas being visited more often. In other cities, the distribution is more uniform, with key intersections being visited the most.

B. Autolabeling

As described in the main paper, We developed an offboard perception system to generate the bounding boxes and tracks for the objects in the scene. It consists of three stages: object detection, offline tracking, and global track refinement. The system is deployed on the raw sensor data from nuPlan dataset.

Object detection. For the first stage, we extend the state-of-the-art MVF++ [88], [16] as our lidar-based 3D object detector. We select n past frames ($n = 7$) and compensate the points based on ego-motion, which significantly densifies the point clouds in the scene. We find that there is no improvement in incorporating future frames as well. To

further boost the performance, we make several changes to the architecture. First, different from other works [16], [89], we propose to include the points from four side lidars in addition to the top lidar, so the objects closer to the ego vehicle will not be missed. To use these points properly, we adopt the cylindrical view [89] in MVF++, instead of the perspective view, to mitigate the collision of points from different lidars. Second, we enlarge the model capacity by using a RegNet [70] backbone. The RegNet backbone processes voxelized outputs from MVF++ and generates expressive feature maps for final detection heads. Third, a CenterPoint [5] detection head is applied to predict the bounding boxes. Since the resulting object boxes are often flipped, we perform majority voting using past and future detections and adjust the heading accordingly.

Offline tracking. Given the instantaneous bounding boxes generated from the first stage, a Kalman filter based multi-object tracker [90] is applied to associate the boxes across timestamps. Since the tracker is running in offline mode, we leverage more information from both past and future frames to manage tracks: First, after Bir_{min} frames of detections are used to confirm a track, the track is initiated from the first frame of Bir_{min} frames, instead of the last one [90]. As a result, the number of False Negatives is reduced. Second, compared with online tracking [90], we increase Age_{max} to maintain a longer memory before coasting a track. While a higher Age_{max} slows down the method, the number of ID switches is reduced significantly.

Global track refinement. Even though the detection network takes n frames of lidar point clouds as input, it still perceives the scene in a short time window. The resulting points are limited to a certain perspective, posing challenges to estimating the size and heading accurately. Moreover, although the bounding boxes from the detection network have been smoothed to some extent by a Kalman Filter based tracker, the smoothing does not happen at a global scale. To mitigate this effect, we introduce a novel neural network. Inspired by recent works [16], [91], our network loads tracks along with the aggregated point clouds within the tracks and refines the bounding boxes in the tracks in terms of position, heading, size and velocity. Notably, instead of using two networks for static and dynamic objects [16], or for different purposes (size, position or heading) [91], we design a single network, which achieves the above goals in one shot. In practice, the proposed method reduces the deployment overhead compared to other works [16], [91], because the order of the two networks is not clear and the classification of dynamic/static is arbitrary for some slow-moving objects. Similar to these works, we also only apply this network to the vehicle class.

Implementation Details. For simplicity, we consider objects on a 2D bird’s eye view plane only and denote (x, y) as the object’s position, (w, l, h) as size, (v_x, v_y) as velocity and θ as heading. An object at timestamp t in a track is denoted as \mathcal{O}_t . It has a bounding box $\mathcal{B}_t = (x_t, y_t, z_t, w_t, h_t, l_t, \theta_t)$ that tightly contains the points \mathcal{P}_t . We select the boxes from the past and future n_b frames. These $(2n_b + 1)$ bounding boxes

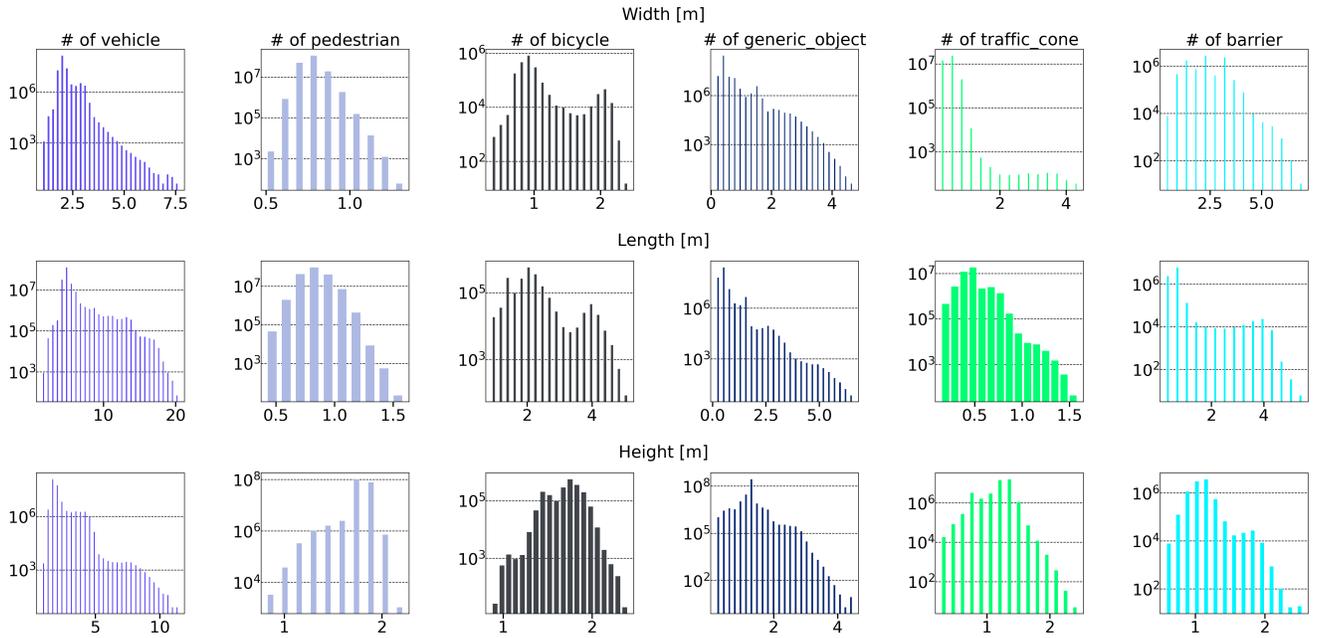


Fig. 7. Size distributions of boxes.

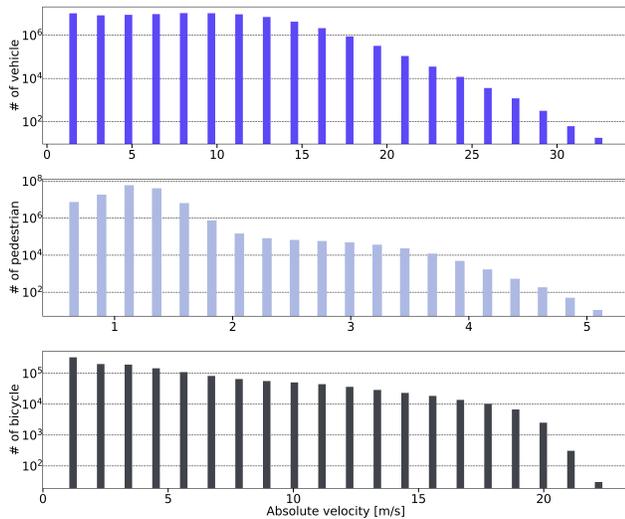


Fig. 8. Histogram of the absolute box velocities per class. Only boxes with speed $> 0.5m/s$ are shown. Note the logarithmic scaling on the y axis.

are concatenated as the input for the trajectory encoding branch. Meanwhile, the point clouds are cropped in both past and future n_p frames and then transformed to the current frame t . The aggregated point cloud is processed by dynamic voxelization [88] and taken as input by the point encoding branch. Convolutional layers and global average pooling are applied for each branch to generate expressive features, which are concatenated for the final MLP. In training, we use a smooth L1 loss to regress the residuals between the ground truth and input \mathcal{B}_t . In deployment, after getting the refined bounding boxes for the whole track, we choose the median

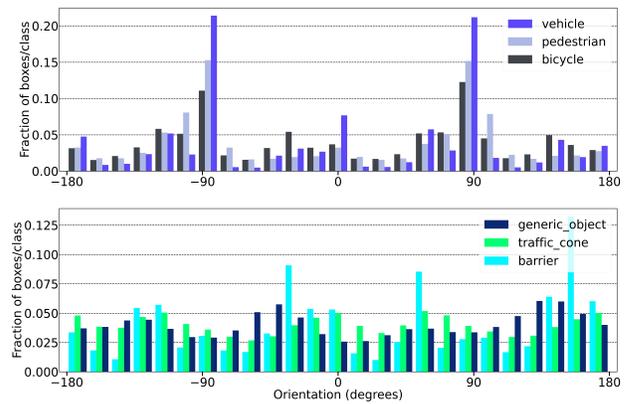


Fig. 9. Orientation of boxes.

value of the size for all the bounding boxes. We update the position and velocity of each box according to the outputs of the network.

C. Evaluation of the autolabeling system

Dataset. To evaluate the proposed autolabeling system we use an internal dataset that follows a similar data distribution to nuPlan and contains 3098 human-labeled scenes collected from Singapore, Boston, Pittsburgh and Las Vegas. Among them, 2958 scenes are used for training, 49 for validation and 91 for testing.

Metrics. The performance of the detection network and the global track refinement are evaluated using the maximum F1 scores for all the classes. The performance of the offline tracker is evaluated using AMOTA, ID switches, recall and

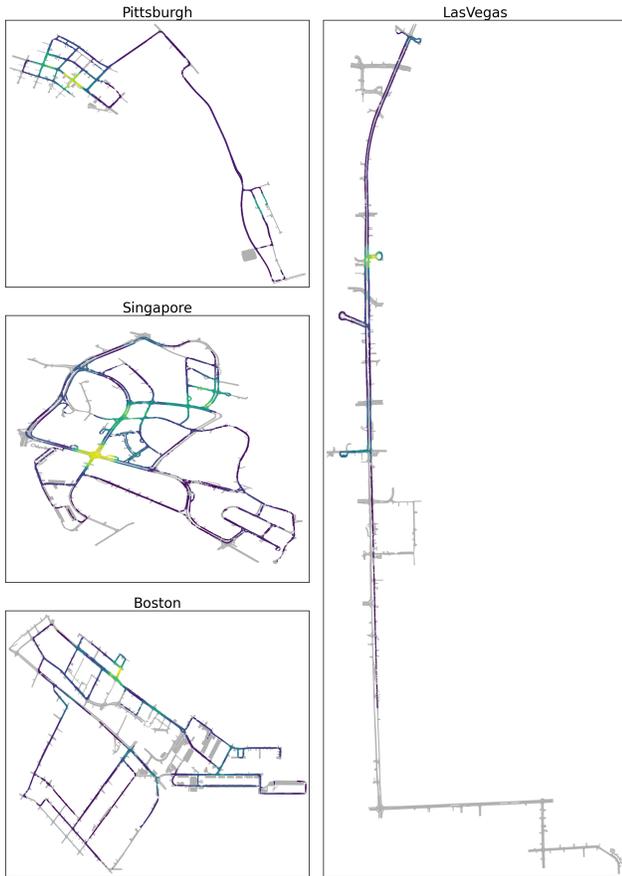


Fig. 10. Spatial data coverage for four cities in nuPlan. Colors indicate the number of scenarios with subsampled ego vehicle poses within a 100m radius. There is more ego vehicles in a location as the colors become more yellow.

F1 score. Vehicles are divided into two categories depending on whether their lengths are larger than 7m. We use Birds Eye View 2D IOU as the matching criterion. The IOU threshold is 0.7 for vehicles, 0.5 for cyclists and pedestrians, 0.3 for barriers and generic objects, 0.15 for traffic cones.

3d object detection. The detection network proposed in this paper is compared with a modified CenterPoint [5] detector, termed *mCenterPoint*. Our goal is to compare our proposed network with a detection network often deployed onboard for autonomous vehicles. We choose *mCenterPoint* as a baseline because it is lightweight and meets the real-time requirement. It differs from our offline detection network in three aspects. First, it aggregates past 3 sweeps of point clouds as input to the network. Second, *mCenterPoint* uses the vanilla PointPillars [2] for voxelization. Third, a more compact backbone is used to extract features for the CenterPoint detection heads. In contrast, the offline model is applied on 7 sweeps of point clouds with a multi-view encoder [16] and a heavy backbone [70]. The max F1 scores for each class are shown in Tab. IX. From the table, we can see the proposed detection network can outperform *mCenterPoint* by a large margin. It clearly shows that using more sweeps of point clouds as well as a more complex network architecture

TABLE IX
DETECTION PERFORMANCE (MAX F1) OF THE PROPOSED AUTOLABELING DETECTOR, COMPARED TO A MODIFIED VERSION OF THE COMMONLY USED CENTERPOINT DETECTOR

Class	mCenterPoint	Ours
Short vehicle	72.0	86.2
Long vehicle	63.9	74.8
Cyclist	64.3	73.0
Pedestrian	59.5	75.1
Traffic cone	68.7	78.3
Barrier	60.0	68.3
Generic object	60.2	68.8
Average	64.1	74.9

TABLE X
COMPARISON OF TRACKING METRICS BETWEEN THE PROPOSED OFFLINE TRACKER AND THE MODIFIED AB3DMOT

Method	AMOTA	ID switch	Recall	F1 score
mAB3DMOT	0.611	9779	0.696	0.828
Ours	0.684	2129	0.727	0.868

improves the performance drastically.

Offline tracker. We also implemented a modified version of AB3DMOT [90] as a baseline in the tracking experiments. Similar to *mCenterPoint*, *mAB3DMOT* is configured for real-time deployment onboard for autonomous vehicles. In particular, Bir_{min} and Age_{max} are set to be 3 and 0.4s for *mAB3DMOT*, respectively. In comparison, the offline tracker uses $Age_{max} = 2s$ and back-traces the first 2 frames with $Bir_{min} = 3$. From Tab. X we see a significant improvement on all tracking metrics. In particular, the number of ID switches has been reduced by 78%, which is due to the significantly enlarged Age_{max} . The fewer number of ID switches indicates that the tracks in nuPlan dataset are less fragmented and thus reflect the agents' movement in the real world.

Global track refinement. To evaluate global track refinement, we compare it against the outputs of the offline tracker. Because the offline tracker will interpolate and suppress some bounding boxes from the detection network, the F1 scores reported here are different from Tab. IX. The results are presented in Tab. XI. It is shown that the vehicle F1 score is consistently improved by global track refinement. In particular, when a stricter matching criterion is applied (BEV IOU = 0.9), we see a relative boost of 56.8%, which demonstrates the efficacy and necessity of global track refinement in producing high-quality bounding boxes.

TABLE XI
COMPARISON OF F1 SCORES OF VEHICLE BOUNDING BOXES BETWEEN THE PROPOSED GLOBAL TRACK REFINEMENT AND THE OFFLINE TRACKER.

Method	IOU 0.7	IOU 0.9
Offline-tracker	83.4	21.3
Global track refinement	87.5	33.4

D. Traffic lights

To automatically label traffic light statuses within nuPlan, we make use of the tracks produced by our autolabeling system, as well as the human-annotated map information. The map indicates the intersections with traffic lights. Within each traffic light intersection, the individual traffic lights control the flow of vehicles from a lane on one side of the intersection to another lane. Each such pair of lanes is connected via a *lane connector* in our map. Hence, we encode the status of the traffic lights in the corresponding lane connectors.

Green and amber statuses. We consider both green and amber statuses of traffic lights to be the same, i.e. green, since vehicles are allowed to move under both light statuses. To infer the presence of a green traffic light on a particular lane, we determine if there are agents moving along the lane connector corresponding to the traffic light. We do this by comparing the directed Hausdorff distance between the trajectories of all the agents within the traffic light intersection, and the lane connector. Note that as the traffic light labeling system is offline, we are able to make use of both the past and future trajectories of each agent. If the directed Hausdorff distance is small, we surmise that there are agents moving along the lane connector, and thus the traffic light controlling that lane connector is likely green.

Red statuses. To infer the presence of a red traffic light on a particular lane, we check for the minimum speed of the agents that are on that lane. Note that we only take into account agents which are of a certain distance from the traffic light intersection. If the minimum speed of the agents on that lane is low, we surmise that the agents are stopped. We also consider the deceleration of the agents on the various lanes. When a traffic light is red, it is common for drivers to begin decelerating as they approach the intersection. Thus, we set a threshold for the deceleration magnitude above which we consider the agent in a lane to be decelerating. If we find that either the agents on a given lane are stopped or decelerating, then we infer that the traffic light controlling that lane connector is likely red. For lanes and lane connectors for which there are no observable agents, we set the status of the corresponding traffic light to ‘unknown’.

Post-processing. After inferring the per-frame green and red statuses, we perform post-processing to refine the inferences. First, we perform grouping. Our map stores information on lane connectors that go in a “parallel” direction, and therefore share the same traffic light statuses. We use this information to set all lane connectors in the parallel direction to have the same status. This helps to reduce false negatives, especially in situations when there are no observed agents on some of the lanes going in the same direction.

Second, we perform back-filling. When a traffic light changes from red to green, drivers often have a certain reaction time before moving off. Similarly, when a traffic light changes from yellow to red, there may be drivers still crossing the intersection. This might lead to the system inferring a green traffic light for the particular frame, even

though the lights have changed in reality. Hence there is a slight lag in the transitions identified via motion inference with respect to the actual transition. To account for this, for all green statuses, we go back a specified time horizon into the past and override the statuses of each lane connector by setting them to green. We do the same for the red statuses and override the past statuses within the specified time horizon with red.

Evaluation. To perform a more quantitative evaluation of our traffic light labeling system, we select scenes where the ego vehicle is moving through traffic light intersections and manually label the traffic light statuses for a subset of the data. We select scenes from various cities and various traffic light behaviors (e.g. constant, transition). We manually label about 1000 frames. We seek to compare our system against a more conventional traffic light detection system. Such a system is usually deep-learning based and vision-only. Consequently, for the baseline, we use YOLOv3 [92], as implemented by [93] and trained on the LISA Traffic Light Dataset [94], which contains 113,888 annotated traffic lights. We do not fine-tune YOLOv3 on any nuPlan data. By labeling the statuses of the traffic lights via motion inferences, we are able to recover 5.2x more traffic light statuses than YOLOv3. We find that YOLOv3 [92] misses several traffic lights at long distances. It also performs poorly at very close distances to the traffic light due to the perspective warping of the cameras. Using YOLOv3 [92], the traffic light classification accuracy among visible traffic lights is 66.4%. In contrast, our system performs slightly better with an accuracy of 68.7%. We find that the transitions of the traffic light statuses are difficult to infer for our system, due to the lag in the agents moving off, slowing down or stopping whenever the traffic lights change. A data-driven alternative could provide better results here.

E. Scenario mining

We propose the following approach to mine scenarios from the nuPlan dataset. First, we compute a large number of atomic primitives. These primitives model an attribute (vehicle speed) or state transition (a vehicle being in two lanes simultaneously) and can be extracted from the entire dataset in a single pass. Second, we combine several primitives into an SQL query that can be run efficiently on a database. Third, we post-process the query results by including a sequence of 10 seconds before and after the returned time step. Finally, we manually QA 100 examples for each scenario. If the false positive rate is below 90%, we either refine the query by adding more attributes or tuning hyperparameters or discard the scenario. This approach results in high-precision scenario labels, while recall is less relevant for us. We have a total of 73 unique scenario types in the dataset, and their distributions per city are shown in Fig. 11. The details and parameters of the 14 scenarios that were used to grade the nuPlan planning challenge can be found in Tab. XII.

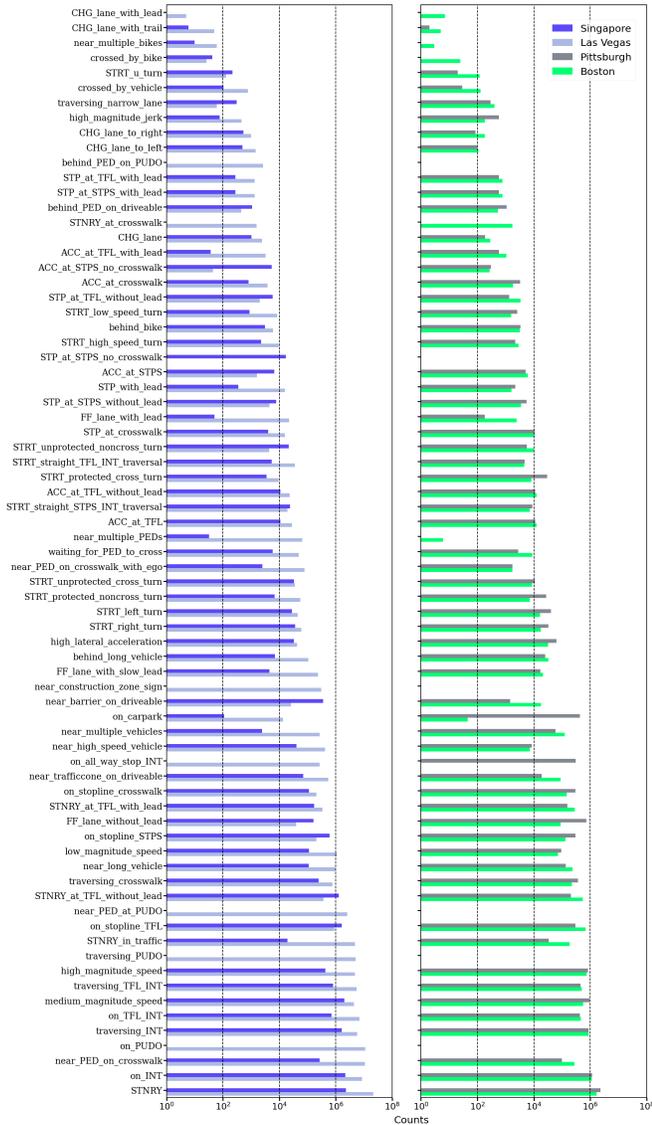


Fig. 11. Distribution of scenario types in each city.

VIII. SIMULATION

A. Agents

The differential equation behind the IDM policy operates based on a focus agent - the vehicle that the IDM policy is controlling - and a lead object. The policy is parameterized by the distance to a lead object. However, in closed-loop simulation identifying the correct lead object is not trivial. Fig. 12 shows how the closest leading object is identified. Searching for the closest object often returns another agent in the adjacent lane. Hence, the search space is reduced to only the other objects in the scene that are or can potentially intersect with the agent's planned path. The euclidean distance between the two closest points of the focus agent and the lead object is computed.

Merging situations can be tricky because IDM does not inherently account for multi-lane interactions. One way to simulate this is to extend all other objects' footprints in the

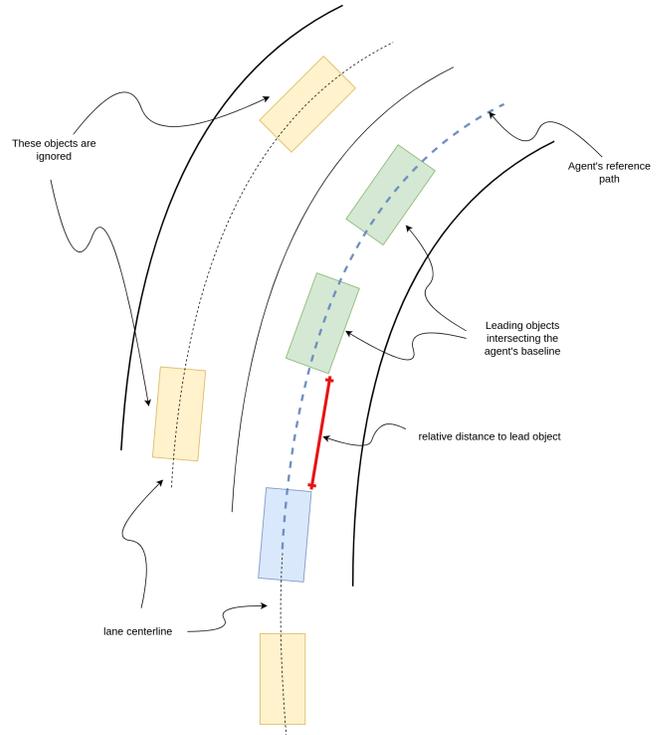


Fig. 12. Filtering for all other objects (green) that intersect with the agent's (blue) baseline. The baseline is expanded into a polygon along the path and as wide as the agent's width. All other objects (yellow) are ignored. Checking for the nearest object out of the filtered ones.

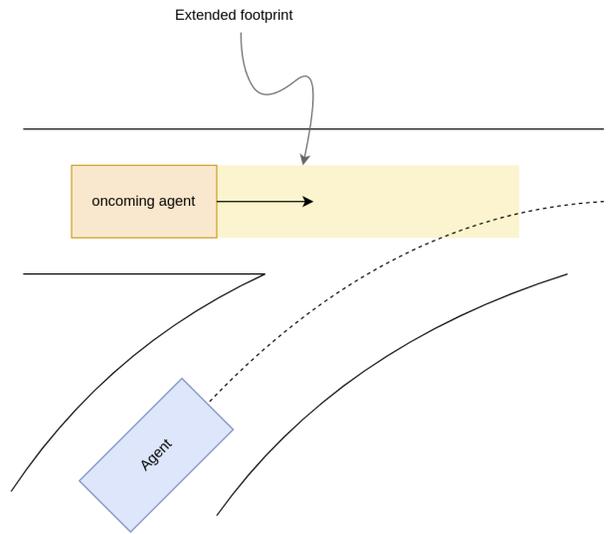


Fig. 13. Footprint projection of other objects in the scene.

TABLE XII
NUPLAN PLANNING CHALLENGE SCENARIO DETAILS

Scenario type	Scenario description
starting straight traffic light intersection traversal	Ego at the start of a traversal going straight across an intersection area controlled by traffic lights while not stopped.
high lateral acceleration	Ego high ego acceleration ($1.5 < acceleration < 3m/s^2$) across the lateral axis with high yaw rate while not turning.
changing lane	Ego at the start of a lane change towards an adjacent lane.
high magnitude speed	Ego high velocity magnitude with low acceleration ($velocity > 9m/s$).
low magnitude speed	Ego low ego velocity magnitude ($0.3 < velocity < 1.2m/s$) with low acceleration while not stopped.
starting left turn	Ego at the start of a traversal turning left across an intersection area while not stopped.
starting right turn	Ego at the start of a traversal turning right across an intersection area while not stopped.
stopping with lead	Ego starting to decelerate ($acceleration\ magnitude < -0.6m/s^2, velocity\ magnitude < 0.3m/s$) with a leading vehicle ahead ($distance < 6m$) at any area.
following lane with lead	Ego following ($velocity > 3.5m/s$) its current lane with a moving leading vehicle ahead on the same lane ($velocity > 3.5m/s, longitudinal\ distance < 7.5m$).
near multiple vehicles	Ego nearby ($distance < 8m$) of multiple (> 6) moving vehicles while ego is moving ($velocity > 6m/s$).
traversing pickup dropoff	Ego during the traversal of a pickup/drop-off area while not stopped.
behind long vehicle	Ego behind ($3m < longitudinal\ distance < 10m$) a long ($length > 8m$) vehicle in the same lane as ego ($lateral\ distance < 0.5m$).
waiting for pedestrian to cross	Ego waiting for a nearby ($distance < 8m, time\ to\ intersection < 1.5m$) pedestrian to cross a crosswalk area while ego is not stopped and the pedestrian is not at a pickup/drop-off area.
stationary in traffic	Ego is stationary with multiple (> 6) vehicles nearby ($distance < 8m$).

scene proportional to their speed. Fig. 13 shows an example of a situation where an extended footprint can help in lane-merge situations. The lead agent search will identify the oncoming agent. The agent will know to slow down sooner, giving way to the oncoming objects.

The final IDM algorithm looks as such:

- 1) Breadth-first search path planning for a set distance.
- 2) Identifying the closest leading object.
- 3) Perform a one-step forward euler numerical integration on the IDM differential equations.
- 4) Propagate the agent along the planned path according to the solution of step 3.
- 5) Repeat for all agents in the scene.
- 6) Repeat for all simulation propagation steps.

The attempt to re-purpose prevailing motion forecasting models such as LaneGCN for traffic simulation proved less trivial than initially thought. The model independently predicts agents, resulting in a lack of scene cohesion, for example, predicting colliding trajectories. The model is also

susceptible to distribution shift issues that arise in closed-loop simulation. When applied to traffic simulation, the model induces unrealistic driving scenarios [95]. Furthermore, the number of simulated agents had to be limited to maintain acceptable simulation runtime. It can be concluded that motion forecasting alone cannot achieve realistic, scene-coherent traffic simulation.

B. Evaluation

Metrics. As explained in the main paper, we design a set of metrics along with a scoring function to compare the performance of planners. The previously mentioned open-loop metrics are described in detail in Tab. XIII. For each metric, we evaluate an aggregated error/MR considering different time horizons within the planning horizon (i.e., $H = 3, 5, 8s$) and with the same sampling frequency of $1Hz$. This method allows us to have a fair comparison across planners with different planning horizons or sampling rates. Additionally, by taking the mean across the selected horizons, the errors at the beginning of the horizon will have more impact on the averaged value. The "within bound" metric score used in the cost structure is found by comparing the average error/MR value to a maximum acceptable threshold ($8m$ for distance errors, $0.8rad$ for heading errors, and 0.3 for MR). It is 0 if the average value is more than the threshold, and 1 otherwise. For example, 'MR within bound' from open-loop metrics, and 'drivable area compliance' from the closed-loop metrics. Multiplier metrics are assigned a score of 1 or 0, except for 'no at-fault collision' which takes 0 (if there is an at-fault collision with a vehicle, bicycle, or pedestrian, or there are multiple at-fault collisions with objects), 0.5 (if there's an at-fault collision with a single object), and 1 (if there's no at-fault collision).

In the following, we include additional information about closed-loop metrics mentioned in the main paper:

- When identifying at-fault collisions, we only penalize the planner when the ego vehicle could be responsible for the collision, which includes collisions with stopped agents, collisions with agents in front of ego and collisions with agents in adjacent lanes while making a lane change. On the other hand, the ego is not penalized for rear-end collisions or other agents colliding with the ego when it is stopped. To further emphasize the importance of different agent types, at-fault collisions are grouped into vulnerable road users (including pedestrians and bicyclists), vehicles and objects (traffic cones, barriers and generic objects).
- For drivable area violation, we measure the maximum distance of the corners of the ego bounding box from the nearest drivable area.
- For driving direction, the movement of the ego during a 1s time horizon is calculated along the driving direction of its lane.
- Speed limit violation is defined based on the magnitude and duration of the violation.
- Time to collision is defined as the time required for ego and another object to collide if they continue at

TABLE XIII
OPEN-LOOP METRICS AND THEIR DEFINITIONS

Metric name	Metric definition
ADE Within Bound	At each sampled time, ADE is defined as the average of pointwise L2 distances between the planner trajectory (x-y) and expert trajectory, up to the selected comparison horizon in the future.
FDE Within Bound	At each sampled time, FDE is defined as the L2 distance between the planner trajectory (x-y) and expert trajectory at the final time available in the sampled trajectories.
AHE Within Bound	At each sampled time, AHE is defined as the average of absolute differences between the planner trajectory heading and expert trajectory heading up to the selected comparison horizon in the future.
FHE Within Bound	At each sampled time, FHE is defined as the absolute differences between the planner trajectory heading and expert trajectory heading at the final time available in the sampled trajectories.
Miss Rate Within Bound	At each sampled time, if the maximum of the pointwise L2 distances between the planner trajectory and expert trajectory up to the selected comparison horizon in the future is greater than its corresponding maximum displacement threshold (6, 8 and 16m for horizons 3, 5 and 8s, respectively), we consider the planner trajectory at that time as a miss. Miss rate is the ratio of sampled times where the trajectory was marked as a miss over the number of the sampled times.

their present speed and heading. We only compute time to collision for objects in front of the ego, cross-traffic objects and lateral objects on the sides, when the ego is making a lane change or is in the intersection.

- Rider comfort is measured based on jerk, acceleration and steering rate which are compared to those observed in human driving.
- Progress of the planner trajectory towards the goal is evaluated by comparing its progress along expert’s route in the same scenario. The metric quantifies the progress as the ratio of overall ego progress to the overall expert’s progress during the scenario.

Final score structure. The final score of a planner is computed by averaging its scores across all scenarios as defined in the main paper. What follows helps explain the equation: For open-loop planners, the driven trajectory in a scenario is assigned a zero score if the miss rate is above the selected threshold (0.3), otherwise, a weighted average of other metrics’ scores is used as the score. All weights were tuned with the objective to maximize the overall performance of the planner against human-driven future trajectories across multiple scenarios. The weights can be found in the main paper. For closed-loop planners, the driven trajectory in a scenario is assigned:

- A zero score if 1) there is an at-fault collision with a vehicle or a VRU, or 2) there are multiple at-fault collisions with objects (e.g. a cone), or 3) there is a drivable area violation, or 4) ego drives into oncoming traffic more than 6m (driving distance), or 5) ego progress towards the destination is smaller than a threshold.
- The weighted average of other metrics’ scores is mul-

TABLE XIV
CLOSED-LOOP METRICS, SCORES AND WEIGHTS

Metric name	Metric score
No at-fault Collisions	0, 0.5 or 1
Drivable Area Compliance	boolean
Driving Direction Compliance	0, 0.5 or 1
Making progress	boolean
Speed limit compliance	$f_{speed-limit}$
Time to Collision within bound	boolean
Progress along route	$\frac{ego\ progress}{expert\ progress}$
Comfort	boolean

tiplied with 0.5 if there is one at-fault collision with an object (e.g. a cone), or if ego drives into oncoming traffic more than 2m, but less than 6m.

- A weighted average of other metrics’ scores, otherwise.

Closed-loop metric scores are summarized in Tab. XIV. $f_{speed-limit}$ is a function that returns 1 if there are no speed limit violations and approaches 0 as the violation increases. Furthermore, the comfort metric accounts for jerk amplitude, lateral and longitudinal acceleration, and jerk, and yaw rate and acceleration. It is assigned a zero score if one of the comfort bounds is violated. Our scoring heuristic described above is an initial proposal that accounts for the natural importance of each metric and is hand-tuned for our dataset and simulation framework.

IX. EXPERIMENTS

A. Planning experiments

Detailed scenario-stratified metrics for all four planner baselines (rule-based and learned) across the three challenges (open-loop, closed-loop non-reactive, and closed-loop reactive) can be found in Tab. XV. The metric breakdown corroborates the statement in the main paper that metrics reward conservative driving. Planners that do not collide and stay within the drivable area may score better than planners that attempt to mimic human drivers. Hence, finer grain and scenario-based metrics are required to distinguish between simple rules-abiding driving from desirable human-like driving behaviors.

B. nuPlan Challenge

The distribution of planners’ scores across the nuPlan challenge can be seen in Fig. 14 for open-loop, Fig. 15 for closed-loop non-reactive and Fig. 16 for closed-loop reactive. Most submitted planners were able to score well in the open-loop challenge. The most common scores lie between 0.8 - 0.85. The scores significantly dropped for closed-loop challenges. The most common scores lie between 0.65 - 0.7. A 0.15 drop from the open-loop challenge. This further supports the argument that most purely learned models fail to generalize to closed-loop scenarios. For an overview of the leaderboard, please refer to ¹.

¹<https://eval.ai/web/challenges/challenge-page/1856/leaderboard/4360>

TABLE XV

THE METRICS BREAKDOWN BY SCENARIO TYPES ACROSS THE THREE CHALLENGES FOR ALL PLANNERS: SIMPLEPLANNER (SP), IDMPLANNER (IDM), RASTER PLANNER (RP), AND URBANDRIVER (UD)

Scenario type	Open-loop				Closed-loop Non-reactive				Closed-loop Reactive			
	SP	IDM	RP	UD	SP	IDM	RP	UD	SP	IDM	RP	UD
all	0.22	0.30	0.52	0.90	0.32	0.73	0.47	0.68	0.37	0.76	0.46	0.67
behind long vehicle	0.74	0.54	0.79	0.95	0.17	0.99	0.79	0.99	0.44	1.00	0.90	0.98
changing lane	0.16	0.29	0.41	0.89	0.35	0.61	0.43	0.81	0.41	0.59	0.42	0.74
following lane with lead	0.15	0.05	0.30	0.89	0.40	0.77	0.48	0.77	0.47	0.91	0.26	0.82
high lateral acceleration	0.07	0.32	0.41	0.86	0.13	0.84	0.24	0.58	0.14	0.85	0.22	0.54
high magnitude speed	0.02	0.15	0.50	0.90	0.64	0.77	0.57	0.89	0.71	0.89	0.47	0.90
low magnitude speed	0.49	0.45	0.66	0.92	0.29	0.83	0.57	0.68	0.38	0.80	0.61	0.79
near multiple vehicles	0.20	0.26	0.37	0.93	0.45	0.82	0.61	0.89	0.52	0.84	0.50	0.80
starting left turn	0.00	0.13	0.39	0.85	0.21	0.56	0.07	0.48	0.21	0.63	0.17	0.51
starting right turn	0.00	0.18	0.41	0.87	0.09	0.33	0.20	0.14	0.09	0.40	0.27	0.12
intersection traversal	0.02	0.21	0.47	0.91	0.55	0.81	0.53	0.78	0.57	0.78	0.44	0.77
stationary in traffic	0.73	0.70	0.87	0.97	0.71	0.97	0.86	0.94	0.73	0.97	0.87	0.93
stopping with lead	0.54	0.54	0.79	0.95	0.09	0.95	0.77	0.91	0.25	0.95	0.72	0.84
traversing pickup dropoff	0.18	0.26	0.45	0.87	0.17	0.68	0.35	0.49	0.17	0.71	0.37	0.48
waiting for pedestrian to cross	0.13	0.17	0.51	0.81	0.09	0.30	0.29	0.22	0.13	0.38	0.37	0.24



Fig. 14. nuPlan challenge open-loop score distribution.

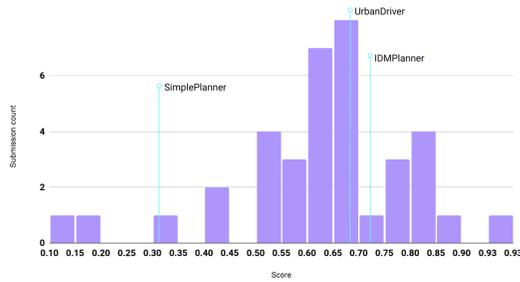


Fig. 15. nuPlan challenge closed-loop non-reactive score distribution.

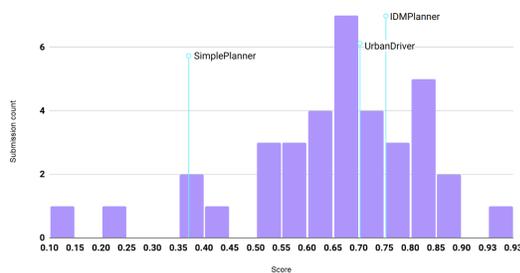


Fig. 16. nuPlan challenge closed-loop reactive score distribution.

TABLE XVI
OPEN-LOOP CHALLENGE

Metric	CS Tu	AutoHorizon	Pegasus	AID
Overall score	0.829	0.852	0.876	0.840
ADE within bound	0.815	0.836	0.843	0.811
FDE within bound	0.597	0.665	0.727	0.639
Miss rate within bound	0.962	0.960	0.960	0.966
AHE within bound	0.944	0.958	0.947	0.937
FHE within bound	0.921	0.938	0.935	0.925

TABLE XVII
CLOSED-LOOP NON-REACTIVE CHALLENGE

Metric	CS Tu	AutoHorizon	Pegasus	AID
Overall score	0.928	0.890	0.817	0.809
Making progress	0.994	0.978	0.929	0.936
Drivable area compliance	1.000	0.990	0.948	0.962
Driving direction compliance	1.000	0.988	0.955	0.992
Comfort	0.919	0.990	0.927	0.940
No at-fault collisions	0.988	0.963	0.926	0.939
TTC within bound	0.925	0.905	0.879	0.883
Progress along expert route	0.914	0.915	0.793	0.841
Speed limit compliance	0.997	0.960	0.934	0.974

TABLE XVIII
CLOSED-LOOP REACTIVE CHALLENGE

Metric	CS Tu	AutoHorizon	Pegasus	AID
Overall score	0.929	0.881	0.851	0.838
Making progress	0.992	0.980	0.946	0.946
Drivable area compliance	1.000	0.992	0.952	0.966
Driving direction compliance	1.000	0.992	0.962	0.996
ego is comfortable	0.925	0.992	0.952	0.972
No at-fault collisions	0.993	0.965	0.946	0.969
TTC within bound	0.954	0.921	0.907	0.919
Progress along expert route	0.885	0.882	0.799	0.828
Speed limit compliance	0.997	0.964	0.940	0.981