

Argoverse: 3D Tracking and Forecasting with Rich Maps

Ming-Fang Chang^{* 1,2}, John Lambert^{*1,3}, Patsorn Sangkloy^{*1,3}, Jagjeet Singh^{*1}, Sławomir Bak¹, Andrew Hartnett¹, De Wang¹, Peter Carr¹, Simon Lucey^{1,2}, Deva Ramanan^{1,2}, and James Hays^{1,3}

¹Argo AI, ²Carnegie Mellon University, ³Georgia Institute of Technology

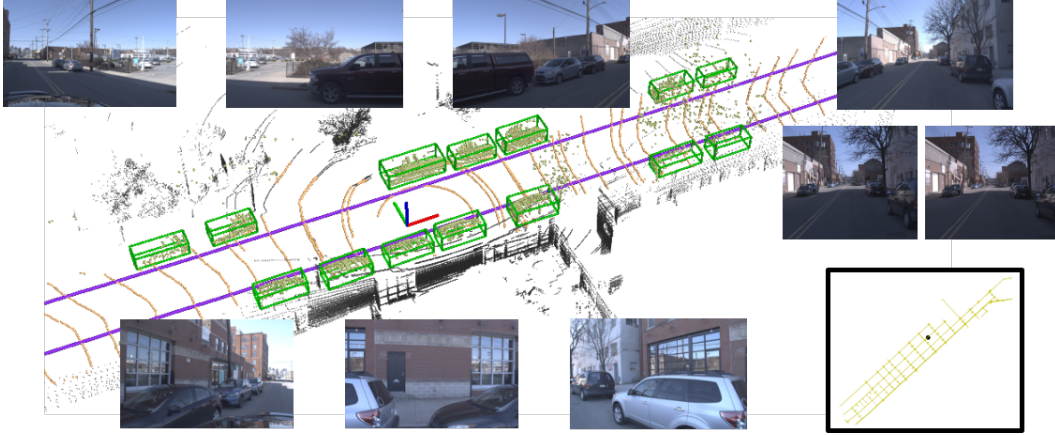


Figure 1: We introduce a dataset for 3D tracking and forecasting with *rich maps* for autonomous driving. Our 3D tracking dataset contains sequences of LiDAR measurements, 360° RGB video, front-facing stereo (middle-right), and 6-dof localization. All sequences are aligned with maps containing lane center lines (magenta), driveable region (orange), and ground height. Sequences are annotated with 3D cuboid tracks (green). A wider map view is shown in the bottom-right.

Abstract

We present Argoverse, a dataset designed to support autonomous vehicle perception tasks including 3D tracking and motion forecasting. Argoverse includes sensor data collected by a fleet of autonomous vehicles in Pittsburgh and Miami as well as 3D tracking annotations, 300k extracted interesting vehicle trajectories, and rich semantic maps. The sensor data consists of 360° images from 7 cameras with overlapping fields of view, forward-facing stereo imagery, 3D point clouds from long range LiDAR, and 6-dof pose. Our 290km of mapped lanes contain rich geometric and semantic metadata which are not currently available in any public dataset. All data is released under a Creative Commons license at www.argoverse.org. In baseline experiments, we use map information such as lane direction, driveable area, and ground height to improve the accuracy of 3D object tracking. We use 3D object tracking to “mine” for more than 300k interesting vehicle trajectories to create a trajectory forecasting benchmark. Motion forecasting experiments ranging in complexity from classical methods (k -NN) to LSTMs demonstrate that using detailed “vector maps” with lane-level information substantially reduces

prediction error. Our tracking and forecasting experiments represent only a superficial exploration of the potential of rich maps in robotic perception. We hope that Argoverse will enable the research community to explore these problems in greater depth.

1. Introduction

Datasets and benchmarks for a variety of perception tasks in autonomous driving have been hugely influential to the computer vision community over the last few years. We are particularly inspired by the impact KITTI [12] has had in opening new research directions. However, publicly available datasets for autonomous driving rarely include *map* data, even though detailed maps are critical to the development real world autonomous systems. Publicly available maps, e.g. OpenStreetMap, can be useful, but have limited detail and accuracy.

Intuitively, 3D scene understanding would be easier if maps directly told us which 3D points belong to the road, which belong to static buildings, what lane a tracked object is in, what the speed limit for that lane is, how far it is to the next intersection, etc. But since publicly available datasets don’t contain such rich mapped attributes it is an

^{*}Equal contribution

open research question of how to represent and utilize these features. Argoverse is the first large autonomous driving dataset with such detailed maps. We examine the potential utility of these new map features on two tasks – 3D tracking and motion forecasting, and we offer a significant amount of real-world, annotated data to enable new benchmarks for these problems.

Our contributions in this paper include:

- We release a large scale dataset with synchronized data from LiDAR, 360° and stereo cameras sampled across two cities and varied conditions.
- We provide ground truth tracking annotation of objects in 3D, with ten times more tracks than the KITTI [12] tracking benchmark.
- We create a large scale forecasting benchmark of trajectories capturing scenarios like turns at intersections, driving with many vehicles nearby, and lane changes.
- We release map data and an API which can be used to develop map-based perception algorithms. To our knowledge, there is no publicly available equivalent to our semantic vector map of road infrastructure and traffic rules.
- We examine the influence of map context in 3D tracking and trajectory forecasting.
- We release the first large-scale dataset suitable for training and benchmarking automatic map creation, often known as *map automation*.
- We release the first fully panoramic, high-frame rate large-scale dataset collected outdoors on a vehicle, opening new possibilities for city-scale reconstruction with photometric-based direct methods.

2. Related Work

Autonomous Driving Datasets with Map Information.

Until recently, it was rare to find datasets that provide detailed map information associated with annotated data. Works such as TorontoCity [50] and ApolloScape [22] focus on map construction tasks but without 3D annotation for dynamic objects. The nuScenes dataset [5] contains maps in the form of binary, rasterized, top-down indicators of region of interest (where region of interest is the union of driveable area and sidewalk). This map information is provided for 1000 annotated vehicle log segments (or “scenes”) in Singapore and Boston. Like nuScenes, Argoverse includes maps of driveable area, but we also include ground height and a “vector map” of lane centerlines and their connectivity.

Autonomous Driving Datasets with 3D Tracking Annotations. Many existing datasets for object tracking focus on pedestrian tracking from image/video sequences [42, 36, 2]. Several datasets provide raw data from self driving car sensors, but without any object annotations [35, 39, 43]. The

ApolloCar3D dataset [47] is oriented toward 3D semantic object keypoint detection instead of tracking. KITTI [12] and H3D [41] offers 3D bounding box and track annotations but does not provide a map and the camera field of view is frontal, rather than 360°. VIPER [44] provides data from a simulated world with 3D track annotations. nuScenes [5] currently provides 360° data and a benchmark for 3D object detection, with tracking annotation also available. The Argoverse-Tracking dataset contains 360° track annotations in 3D space aligned with detailed map information. See Table 1 for a comparison between 3D autonomous vehicle datasets.

Autonomous Driving Datasets with Mined Trajectory Data.

TrafficPredict [34] also uses sensor-equipped vehicles to observe driving trajectories in the wild and build a forecasting benchmark. The TrafficPredict dataset consists of 155 minutes of observations compared to 320 hours of observations in Argoverse.

Using Maps for Self-driving Tasks. While high definition (HD) maps are widely used by motion planning systems, few works explore the use of this strong prior in perception systems [52] despite the fact that the three winning entries of the 2007 DARPA Urban Challenge relied on a DARPA-supplied map – the *Route Network Definition File* (RNDP) [37, 49, 3]. Hecker *et al.* [17] show that end-to-end route planning can be improved by processing rasterized maps from OpenStreetMap and TomTom. Liang *et al.* [30] demonstrate that using road centerlines and intersection polygons from OpenStreetMap can help infer crosswalk location and direction. Yang *et al.* [52] show that incorporating ground height and road segment into LiDAR points can improve 3D object detection. Suraj *et al.* [33] use dashboard-mounted monocular cameras on a fleet of vehicles to build a 3D map via city-scale structure-from-motion for localization of ego-vehicles and trajectory extraction.

3D Object Tracking. In traditional approaches for point cloud tracking, segments of points can be accumulated using clustering algorithms such as DBSCAN [11, 28] or connected components of an occupancy grid [29, 21], and then associated based on some distance function using the Hungarian algorithm. Held *et al.* utilize probabilistic approaches to point cloud segmentation and tracking [18, 20, 19]. Recent work demonstrates how 3D instance segmentation and 3D motion (in the form of 3D scene flow, or per-point velocity vectors) can be estimated directly on point cloud input with deep networks [51, 31]. Our dataset enables 3D tracking with sensor fusion in a 360° frame.

Trajectory Forecasting: Spatial context and social interactions can influence the future path of pedestrians and cars. Social-LSTM[1] proposes a novel pooling layer to capture social interaction of pedestrians. Social-GAN [14] attempts to model the multimodal nature of the predictions. However, both have only been tested on pedestrian trajectories,

DATASET NAME	MAP TYPE	EXTENT OF ANNOTATED LANES	DRIVEABLE AREA COVERAGE	CAMERA FRAME RATE	360° CAMERAS	INCLUDE STEREO	# TRACKED OBJECTS
KITTI [12]	None	0 km	0 m ²	10 Hz	no	✓	917
Oxford RobotCar [35]	None	0 km	0 m ²	11/16Hz	no	no	0
H3D [41]	None	0 km	0 m ²	30 Hz	no	no	13,763
nuScenes v1.0 [5]	Raster	0 km	1,115,844 m ²	12 Hz	✓	no	64,386
ApolloScape Tracking	Vector+Raster	0 km	0 m ²	n/a	no	no	21,234+
Argoverse-Tracking-Beta (human annotated)	Vector +Raster	204 km (MIA) +86 km (PIT)	1,074,614 m ²	30 Hz	✓	✓	10,726
Argoverse-Forecasting (mined trajectories)	Vector +Raster	204 km (MIA) +86 km (PIT)	1,074,614 m ²	-	no	no	16.4M

Table 1: **Public self-driving datasets.** We compare recent, publicly available self-driving datasets with 3D object annotations for tracking. Coverage area for nuScenes is based on its *road and sidewalk* raster map. Argoverse coverage area is based on our *driveable area* raster map.

with no use of static context (e.g. a map). Deo et al. [10] propose a convolutional social pooling approach wherein they first predict the maneuver and then the trajectory conditioned on that maneuver. In the self-driving domain, the use of spatial context is of utmost importance and it can be efficiently leveraged from the maps. Chen et al. [8] use a feature-driven approach for social and spatial context by mapping the input image to a small number affordances of a road/traffic state. However, they limit their experiments to a simulation environment. IntentNet [6] extends the joint detection and prediction approach of Luo et al. [32] by discretizing the prediction space and attempting to predict one of eight common driving maneuvers. DESIRE [27] demonstrates a forecasting model capturing both social interaction and spatial context. The authors note that the benefits from these two additional components are small on the KITTI dataset, attributing this to the minimal inter-vehicle interactions in the data.

3. The Argoverse Dataset

Our sensor data, maps, and annotations are the *primary contribution* of this work. We also develop an API which helps connect the map data with sensor information e.g. ground point removal, nearest centerline queries, and lane graph connectivity; see Supplemental Material for more details. Our data and API are available under a Creative Commons and MIT license, respectively, see www.argoverse.org.

We collect raw data from a fleet of autonomous vehicles in Pittsburgh, Pennsylvania, USA and Miami, Florida, USA. These cities have distinct climate, architecture, infrastructure, and behavior patterns. The captured data spans different seasons, weather conditions, and times of day. The data used for our dataset traverses nearly 300km of mapped road lanes and comes from a subset of our fleet operating area.

Sensors. Our cars are equipped with two roof-mounted

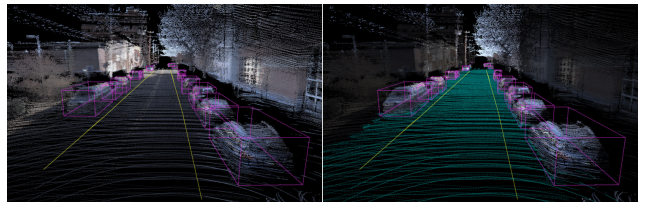


Figure 2: **3D visualization of an Argoverse scene.** Left: we accumulate LiDAR points and project them to a virtual image plane. Right: using our map, LiDAR points beyond driveable area are dimmed and points near the ground are highlighted in cyan. Cuboid object annotations and road center lines are shown in pink and yellow.

VLP-32 LiDAR sensors with an overlapping 40° vertical field of view and a range of 200m, roughly twice that as the sensors used in nuScenes and KITTI. On average, our LiDAR sensors produce a point cloud at each sweep with three times the density of the LiDAR sweeps in the nuScenes [5] dataset (ours $\sim 107,000$ points vs. nuScenes' $\sim 35,000$ points). Each 3D point is motion-compensated to account for ego-vehicle motion throughout the duration of the sweep capture. The vehicles have 7 high-resolution ring cameras (1920×1200) recording at 30 Hz with overlapped field of view providing 360° coverage. In addition there are 2 front-facing stereo cameras (2056×2464) sampled at 5 Hz. Faces and license plates are procedurally blurred in camera data to maintain privacy. Finally, 6-DOF localization for each timestamp comes from a combination of GPS-based and sensor-based localization. Vehicle localization and maps use a city-specific coordinate system described in more detail in the Supplemental Material. Sensor measurements for particular driving sessions are stored in “logs”, and we provide intrinsic and extrinsic calibration data for the LiDAR sensors and all 9 cameras for each log. Figure 2 visualizes our sensor data in 3D. Similar to [43], we place the origin

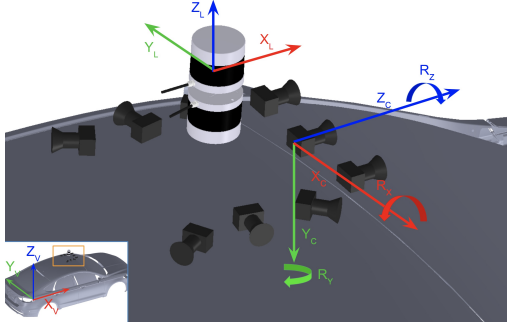


Figure 3: **Car sensor schematic.** Three reference coordinate systems are displayed: (1) the *vehicle frame*, with X_v forward, Y_v left, and Z_v up, (2) the *camera frame*, with X_c across imager, Y_c down imager, and Z_c along optical axis. (3) the *LiDAR frame*, with X_L forward, Y_L left, and Z_L up. Positive rotations R_X , R_Y , R_Z are defined for each coordinate system as rotation about the respective axis following the right-hand rule.

of the vehicle coordinate system at the center of the rear axle. All sensors are roof-mounted, with a LiDAR sensor surrounded by 7 “ring” cameras (clockwise: facing front center, front right, side right, rear right, rear left, side left, and front left) and 2 stereo cameras. Figure 3 visualizes the geometric arrangement of our sensors.

3.1. Maps

Argoverse contains three distinct map components – (1) a vector map of lane centerlines and their attributes, (2) a rasterized map of ground height, and (3) a rasterized map of driveable area and region of interest (ROI).

Vector Map of Lane Geometry. Our *vector map* consists of semantic road data represented as a localized graph rather than rasterized into discrete samples. The vector map we release is a simplification of the map used in fleet operations. In our vector map, we offer lane centerlines, split into lane segments. We observe that vehicle trajectories generally follow the center of a lane so this is a useful prior for tracking and forecasting.

A lane segment is a segment of road where cars drive in single-file fashion in a single direction. Multiple lane segments may occupy the same physical space (e.g. in an intersection). Turning lanes which allow traffic to flow in either direction would be represented by two different lanes that occupy the same physical space.

For each lane centerline, we provide a number of semantic attributes. These lane attributes describe whether a lane is located within an intersection or has an associated traffic control measure (Boolean values that are not mutually inclusive). Other semantic attributes include the lane’s turn direction (*left*, *right*, or *none*) and the unique identifiers for the lane’s predecessors (lane segments that come before) and successors (lane segments that come after) of which

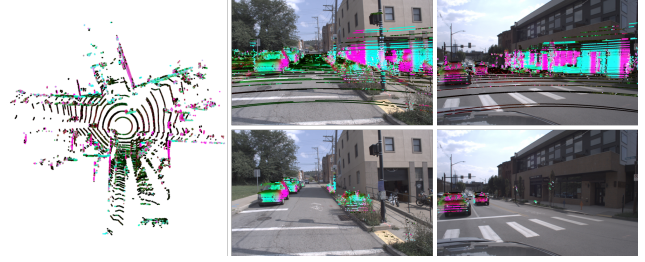


Figure 4: **Map-based ground removal example.** Some Argoverse scenes contain uneven ground, which is challenging to remove with simple heuristics (e.g. assuming ground is planar). Above, the projected LiDAR points are colored by surface normal. The ground surface normal color is non-uniform in the birds-eye-view projection (left). The green color on the slope (middle column) differs from other parts of ground (right column). The lower row uses our map tools to remove ground points and points beyond driveable area.

there can be multiple (for merges and splits, respectively). Centerlines are provided as “polylines”, i.e. a sequence of straight segments. Each straight segment is defined by 2 vertices: (x, y, z) start and (x, y, z) end. Thus, curved lanes are approximated with a set of straight lines.

We observe that in Miami, lane segments that could be used for route planning are on average $3.84\text{m} \pm 0.89$ wide. In Pittsburgh, the average width is $3.97\text{m} \pm 1.04$ in width. Other types of lane segments that would not be suitable for self-driving, e.g. bike lanes, can be as narrow as 0.97m in Miami and as narrow as 1.06m in Pittsburgh.

Rasterized Driveable Area Map. Our maps include binary driveable area labels at 1 meter grid resolution. A driveable area is an area where it is possible for a vehicle to drive (though not necessarily legal). Driveable areas can encompass a road’s shoulder in addition to the normal driveable area that is represented by a lane segment. Our track annotations (Section 3.2) extend to 5 meters beyond the driveable area. We call this larger area our *region of interest* (ROI).

Rasterized Ground Height Map. Finally, our maps include real-valued ground height at 1 meter resolution. Knowledge of ground height can be used to remove LiDAR returns on static ground surfaces and thus makes the 3D detection of dynamic objects easier. Figure 4 demonstrates the use of our ground height map to remove LiDAR points on the road.

3.2. 3D Track Annotations

*Argoverse-Tracking-Beta*¹ contains 113 vehicle log segments with human-annotated 3D tracks. These 113 segments vary in length from 15 to 30 seconds and collectively contain 10,726 tracked objects. We compare this to other datasets in Table 1. For each log segment, we annotate all

¹We refer to our tracking data as *beta* in anticipation of minor refinements or expansion to this dataset before final benchmark release.

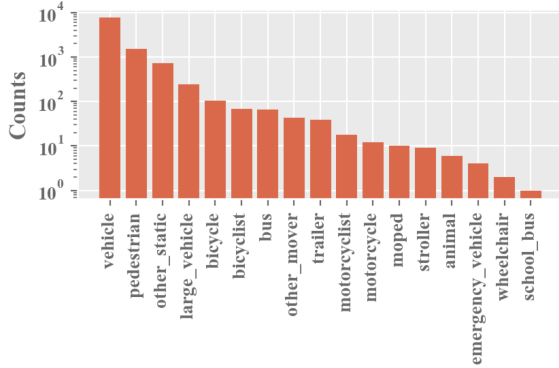


Figure 5: **Distribution of object classes.** This plot shows, in log scale, the number of objects annotated for each class in the 113 log segments in *Argoverse-Tracking-Beta*.

objects of interest (both dynamic and static) with bounding cuboids which follow the 3D LiDAR returns associated with each object over time. We only annotate objects within 5 meters of the *driveable area* as defined by our map. For objects that are not visible for the entire segment duration, tracks are instantiated as soon as the object becomes visible in the LiDAR point cloud and tracks are terminated when the object ceases to be visible. Each object is labeled with one of 17 categories, including ON_ROAD_OBSACLE and OTHER_MOVER for static and dynamic objects that do not fit into other predefined categories. More than 70% of tracked objects are vehicles, but we also observe pedestrians, bicycles, mopeds, and more. Figure 5 show the distribution of classes for annotated objects. All track labels pass through a manual quality assurance review process. Figures 1 and 2 show qualitative examples of our human annotated labels. We divide our annotated tracking data into 65 training, 24 validation, and 24 testing sequences.

3.3. Mined Trajectories for Motion Forecasting

We are also interested in studying the task of *motion forecasting* in which we predict the location of a tracked object some time in the future. Motion forecasts can be critical to safe autonomous vehicle motion planning. While our human-annotated 3D tracks are suitable training and test data for motion forecasting, the motion of many of vehicles is relatively uninteresting – in a given frame, most cars are either parked or traveling at nearly constant velocity. Such tracks are hardly a representation of real forecasting challenges. We would like a benchmark with more diverse scenarios e.g. managing an intersection, slowing for a merging vehicle, accelerating after a turn, stopping for a pedestrian on the road, etc. To sample enough of these *interesting* scenarios we track objects from 1006 driving hours across both Miami and Pittsburgh and find vehicles with interesting behavior in 320 of those hours. In particular, we look for vehicles that are either (1) at intersections (2) taking left or right turns (3) changing to adjacent lanes or (4)

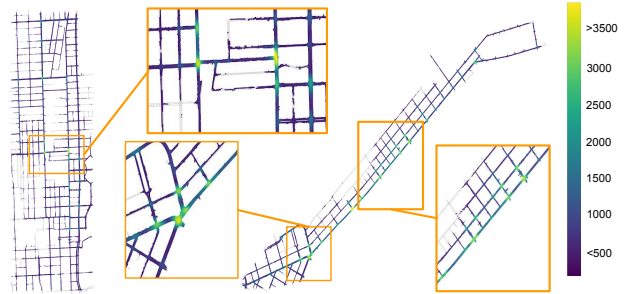


Figure 6: **Distribution of mined trajectories.** The colors indicate the number of mined trajectories across the maps of Miami (left) and Pittsburgh (right). The heuristics to find interesting vehicle behavior lead to higher concentrations in intersections and on busy roads such as Liberty and Penn Ave (southeast roads in bottom right inset).

in dense traffic. In total, we collect 327,793 five second sequences and use them in the forecasting benchmark. Each sequence contains the 2D, birds-eye-view centroid of each tracked object sampled at 10hz. Figure 6 shows the geographic distribution of these sequences. In Section 5, we do not evaluate motion forecasts for pedestrians and stationary vehicles, but still retain their trajectories for context in “social” forecasting models. The 327,793 sequences are split into 208,273 train, 40,128 validation, and 79,392 test sequences. Each sequence has one challenging trajectory which is the focus of our forecasting benchmark. The train, val, and test sequences are taken from disjoint parts of our cities, i.e. roughly one eighth and one quarter of each city is set aside as validation and test data, respectively. This dataset is far larger than what could be mined from publicly available autonomous driving datasets and we have the advantage of using our maps to make it easier to track objects. While data of this scale is appealing because it allows us to see rare behaviors and train complex models, it is too large to exhaustively verify the accuracy of the mined trajectories and thus there is some noise and error inherent in the data.

4. 3D Object Tracking

In this section, we examine how various baseline tracking methods perform on the Argoverse 3D tracking benchmark. Our baseline methods are LiDAR-centric and operate directly in 3D. In addition to measuring the baseline difficulty of our benchmark, we measure how some simple map-based heuristics can influence tracking accuracy. For these baselines, our tracking and evaluation is limited to *vehicles* only.

Given a sequence of F frames, each frame contains set of 3D points from LiDAR $\{P_i \mid i = 1, \dots, N\}$, where $P_i \in R^3$ of x, y, z coordinates, we want to determine a set of track hypothesis $\{T_j \mid j = 1, \dots, n\}$ where n is the

number of unique objects in the whole sequence, and T_j contains the set of object center locations at frames f for $f = \{f_{start}, \dots, f_{end}\}$, the range of frames where the object is visible. We usually have a dynamic observer as our car is in motion more often than not. The tracked vehicles in the scene around us can be static or moving.

Our baseline tracking pipeline clusters LiDAR returns to detect potential objects, uses Mask R-CNN [15] to prune non-vehicle LiDAR returns, associates clusters over time using the Hungarian algorithm, estimating transformations between clusters with ICP, and estimates vehicle pose with a Kalman Filter. More details are provided in the Supplemental Material.

The tracker uses the following map attributes:

Driveable area. Since our baseline is focused on vehicle tracking, we constrain our tracker to the driveable area as specified by the map. This covers any region where it is possible for vehicle to drive (see Section 3.1). This reduces the opportunities for false positives.

Ground removal. We use map information to perform ground-removal. In contrast to local ground-plane estimation methods, the map-based approach is effective in sloping and uneven environments.

Lane Direction. Determining the vehicle orientation from LiDAR alone is a challenging task even for humans due to LiDAR sparsity and partial views. We observe that vehicle orientation rarely violates lane direction, especially so outside of intersections. Fortunately, such information is available in our dataset, so we adjust vehicle orientation based on lane direction whenever the vehicle is not at the intersection and contains too few LiDAR points.

4.1. Evaluation

We use standard evaluation metrics commonly used for multiple object trackers (MOT) [36, 4]. The MOT metric relies on a distance/similarity function between ground truth and predicted objects to determine an optimal assignment. Instead of IoU (Intersection-over-Union) which is more commonly used in tracking literature, we use Euclidean distance between object centroids (threshold for missed track at 2.25 meters, which is half of an average family car length in US). We follow the original definition in CLEAR MOT [4] for MOTP (the lower the better). The tracking metrics are explained in the Supplementary Material in detail.

In the experiments, we run our tracker over the 24 logs in the *Argoverse-Tracking-Beta* test set. We are also interested in the relationship between tracking performance and distance. We apply a threshold (30,50,100 *m*) to the distance between vehicles and our ego vehicle and only evaluate annotations and tracker output within that range. The results in Table 2 show that our tracker performs quite well at short range where the LiDAR sampling density is higher, but struggles for objects beyond 50 meters.

We compare our baseline tracker with three ablations

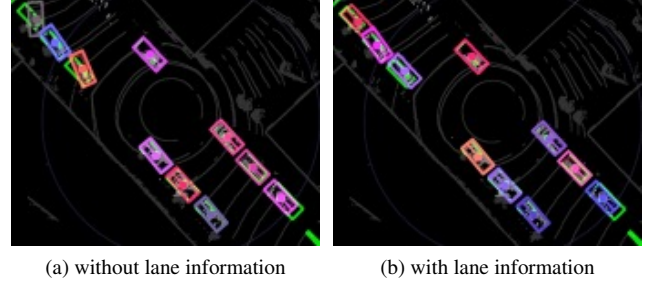


Figure 7: **Tracking with orientation snapping.** Using lane direction information helps to determine the vehicle orientation for detection and tracking. Ground truth cuboids are green.

that exclude: 1) Mask R-CNN as pre-filtering for LiDAR 2) lane direction information from the map and 3) map-based ground removal. The results in Table 2 show that Mask-RCNN dramatically improves our detection performance by reducing false positives. Map-based ground removal leads to slightly better detection performance (higher MOTA) than a plane-fitting approach at longer ranges. On the other hand, lane direction from the map doesn’t affect our metrics (based on centroid distance), but it helps initialize vehicle direction, as shown in Figure 7.

We have used relatively simple baselines to track objects in 3D. We believe that our data opens possibilities in map-based and multimodal tracking research.

5. Forecasting

In this section, we describe our pipeline for trajectory forecasting baselines.

1. Preprocessing: As described in Section 3.3, we first mine for “interesting” sequences and then filter out stationary cars from those. Each sequence contains the centroids of tracked objects over 5 seconds.

Forecasting Coordinate System and Normalization. The coordinate system we use for trajectory forecasting is a top-down, bird eye view (BEV). There are three reference coordinate frames of interest to forecasting: (1) The raw trajectory data is stored and evaluated in the *city* coordinate system (See Section 1.1. of the Supp. Material). (2) For models using lane centerlines as a reference path, we define a *2-d curvilinear coordinate system* with axes tangential and perpendicular to the lane centerline. (3) For models without the reference path (without a map), we align everything such that the observed portion of the trajectory starts at the origin and ends somewhere on the positive x axis. If (x_i^t, y_i^t) represent coordinates of trajectory V_i at timestep t , then this makes sure $y_i^{T_{obs}} = 0$, where T_{obs} is last observed timestep of the trajectory (Section 5.1). We find this normalization works better than leaving trajectories in absolute map coordinates or absolute orientations.

2. Feature Engineering: We define additional features

RANGE THRESHOLD	USE MASK-RCNN	USE MAP LANE	GROUND REMOVAL	MOTA	MOTP	IDF1	MT(%)	ML(%)	# FP	#FN	IDSW	#FRAG
100 m	Y	Y	map	38.01	0.52	0.46	0.10	0.50	103.95	2455.55	32.35	23.20
	N	Y	map	16.45	0.54	0.46	0.16	0.41	1338.95	1972.85	43.70	29.90
	Y	N	map	37.97	0.52	0.46	0.10	0.51	104.90	2455.75	32.10	23.00
	Y	Y	plane-fitting	37.46	0.52	0.46	0.10	0.52	104.40	2483.25	30.65	21.45
50 m	Y	Y	map	52.79	0.52	0.58	0.22	0.29	97.80	1308.75	31.30	22.45
	N	Y	map	21.58	0.54	0.55	0.38	0.18	1196.75	898.45	38.25	24.65
	Y	N	map	52.73	0.52	0.58	0.22	0.29	98.90	1308.85	31.05	22.20
	Y	Y	plane-fitting	52.20	0.52	0.58	0.20	0.31	97.05	1334.85	29.75	20.65
30 m	Y	Y	map	73.12	0.53	0.74	0.66	0.08	91.25	351.55	19.50	13.55
	N	Y	map	23.49	0.57	0.63	0.78	0.04	835.50	238.45	19.50	11.15
	Y	N	map	73.02	0.53	0.74	0.65	0.08	91.90	351.75	19.35	13.45
	Y	Y	plane-fitting	73.14	0.53	0.74	0.66	0.09	89.85	362.35	19.35	13.15

Table 2: **Tracking accuracy at different ranges.** From top to bottom, accuracy for vehicles within 100m, 50m, and 30m.

to capture social and/or spatial context. For social context, we use minimum distance to the objects in front, in back, and the number of neighbors. Such heuristics are meant to capture the social interaction between vehicles. For spatial context, we compute everything in the lane segment coordinate system. We compute the lane centerline corresponding to each trajectory and then map (x_i^t, y_i^t) coordinates to distance along the centerline (a_i^t) and offset from the centerline (o_i^t). In the subsequent sections, we denote social features and map features for trajectory V_i at timestep t by s_i^t and m_i^t , respectively.

3. Prediction Algorithm: We implement weighted Nearest Neighbors and LSTM Encoder-Decoder [40, 13, 48] models using different combinations of features. The results are analyzed in Section 5.3.

5.1. Problem Description

The forecasting task is framed as: *given the past input coordinates of a vehicle trajectory V_i as $X_i = (x_i^t, y_i^t)$ for time steps $t = \{1, \dots, T_{obs}\}$, predict the future coordinates $Y_i = (x_i^t, y_i^t)$ for time steps $\{t = T_{obs}+1, \dots, T_{pred}\}$.* For a car, 5 seconds is sufficient to capture the required part of trajectory, e.g. crossing an intersection. Furthermore, it is unlikely for a typical driving maneuver to last more than 5 seconds. In this paper, we define the forecasting task as observing 20 past frames (2 seconds) and then predicting 10-30 frames (1-3 seconds) into the future. Each trajectory can leverage the trajectories of other vehicles in the same sequence to capture the social context and map information for spatial context.

5.2. Multimodal Evaluation

Predicting the future is difficult. Often, there are several plausible future actions for a given observation. In the case of autonomous vehicles, it is important to predict *many* plausible outcomes and not simply the *most likely* outcome. While some prior works have evaluated forecasting in a deterministic, unimodal way, we believe a better approach is to follow the evaluation methods of DESIRE [27] and Social GAN [14] and encourage algorithms to output multiple predictions.

Our vector map is a semantic graph. The first step in

prediction with a vector map is to localize oneself on the semantic graph. We define two subsequent phases: (1) a *hypothesis phase* and (2) a *generation phase*. The semantic graph makes the generation phase trivial because we can quickly generate hypothesized trajectories via Breadth-First-Search on the semantic graph. However, the hypothesis phase is still challenging due to the multimodal nature of the problem, e.g. it's difficult to know *which* lane segment a vehicle will follow in an intersection.

Among the variety of metrics evaluated in DESIRE was the *oracle error over top K number of samples* metric, where $K = 50$. We follow the same approach and use *top-K Average Displacement Error (ADE)* and *Final Displacement Error (FDE)* as our metrics. The map-based baselines that we report have access to a semantic vector map. As such, they can generate K different hypotheses based on the branching of the road network along a particular observed trajectory. On average, our heuristics generate $K = 5.9$ hypotheses. We generate more than 25 hypotheses for less than 2% of the scenarios. Our map gives us an easy way to produce a compact yet diverse set of forecasts. Other baselines don't have such an option and are restricted to a single prediction. We also provide an *oracle* version of the map-based baselines wherein the model produces the best possible hypothesis by having access to (x_i^t, y_i^t) for $t = \{T_{obs}+1, \dots, T_{pred}\}$, along with the observed trajectory. Note that an oracle-based hypothesis can still generate an imperfect trajectory, e.g. if a car wasn't following any lane.

5.3. Results

In this section, we evaluate the effect of adding social context and spatial context (from the vector map) to improve trajectory forecasting over horizons of 1 and 3 seconds into the future. We evaluate the following models:

- *Constant Velocity:* Compute the mean velocity (v_{xi}, v_{yi}) from $t = \{1, \dots, T_{obs}\}$ and then forecast (x_i^t, y_i^t) for $t = \{T_{obs}+1, \dots, T_{pred}\}$ using (v_{xi}, v_{yi}) as the constant velocity.
- *NN:* Weighted Nearest Neighbor regression where trajectories are queried by (x_i^t, y_i^t) for $t = \{1, \dots, T_{obs}\}$.

BASELINE	1 SECOND		3 SECONDS	
	ADE	FDE	ADE	FDE
Constant Velocity	1.04	1.89	3.55	7.89
NN	0.75	1.28	2.46	5.60
NN+map(oracle)	0.82	1.39	2.39	5.05
NN+map	0.72	1.33	2.28	4.80
LSTM ED	0.68	1.78	2.27	5.19
LSTM ED+social	0.69	1.20	2.29	5.22
LSTM ED+map(oracle)	0.82	1.38	2.32	4.82
LSTM ED+map	0.80	1.35	2.25	4.67
LSTM ED+social+map(oracle)	0.89	1.48	2.46	5.09

Table 3: Forecasting Errors for different prediction horizons

- *NN+map(oracle)*: Weighted Nearest Neighbor regression where trajectories are queried by (a_i^t, o_i^t) for $t = \{1, \dots, T_{obs}\}$ obtained from oracle centerline.
- *NN+map*: Similar to *NN+map(oracle)* but uses *top-K* hypothesized centerlines.
- *LSTM ED*: LSTM Encoder-Decoder model where the input is (x_i^t, y_i^t) for $t = \{1, \dots, T_{obs}\}$ and output is (x_i^t, y_i^t) for $t = \{T_{obs}+1, \dots, T_{pred}\}$
- *LSTM ED+social*: Similar to *LSTM ED* but with input as (x_i^t, y_i^t, s_i^t) , where s_i^t denotes social features
- *LSTM ED+map(oracle)*: Similar to *LSTM ED* but with input as (a_i^t, o_i^t, m_i^t) and output as (a_i^t, o_i^t) , where m_i^t denotes the map features obtained from oracle centerline. Distances (a_i^t, o_i^t) are then mapped to (x_i^t, y_i^t) for evaluation.
- *LSTM ED+map*: Similar to *LSTM ED+map(oracle)* but uses *top-K* hypothesized centerlines.
- *LSTM ED+social+map (oracle)*: Similar to *LSTM ED+map(oracle)* but with input features being $(a_i^t, o_i^t, s_i^t, m_i^t)$.

The results of these baselines are reported in Table 3. Below, we focus on the ADE and FDE for a prediction horizon of 3 seconds to understand which baselines are less impacted by accumulating errors. *Constant Velocity* is outperformed by all the other baselines because it cannot capture typical driving behaviors like acceleration, deceleration, turns etc. *NN+map* has lower ADE and FDE than *NN* because it is leveraging useful cues from the vector map. *NN+map* has lower error than *NN+map(oracle)* as well, emphasizing the multimodal nature of predictions. *LSTM ED* does better than *NN*. *LSTM ED+social* performs similar to *LSTM ED*, implying that the social context does not add significant value to forecasting. A similar observation was made on KITTI [12] in DESIRE [27], wherein their model with social interaction couldn't outperform the one without it. We observe that *LSTM ED+map* outperforms all the other baselines for a prediction horizon of 3 sec. This proves the importance of having a vector map for distant future prediction and making multimodal predictions. Moreover, *NN+map* has a lower FDE than *LSTM ED+social* and *LSTM ED* for higher prediction horizon (3 secs). This sug-

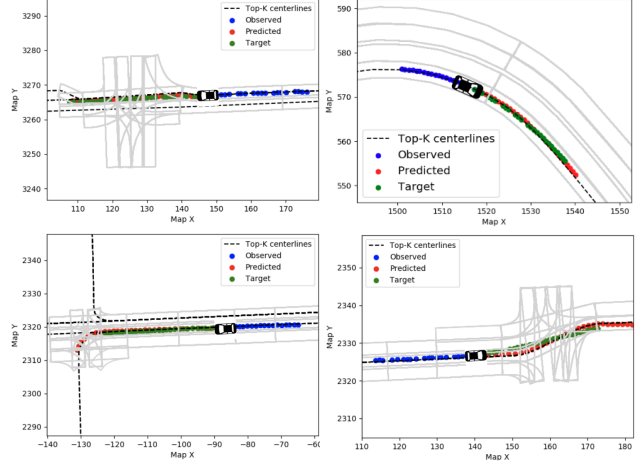


Figure 8: **Qualitative results from *LSTM ED+map* forecasting baseline.** Top left: the model correctly predicts that the car will go straight at the intersection. Top right: the model correctly predicts a smooth right turn never going out of the lane, which might have been difficult if there were no map. Bottom left: demonstration of the multimodal nature of predictions, where the model considers all top-K possibilities. Bottom right: the predictions are on a non-typical lane which takes a slight left and then a slight right. Again, this is hard to predict without a map.

gests that even a shallow model working on top of a vector map works better than a *deep* model with social features and no vector map. Figure 8 shows qualitative forecasting results from our best performing model.

6. Discussion

Argoverse is a large dataset for autonomous driving research. Unique among such datasets, Argoverse contains rich map information such as lane centerlines, ground height, and driveable area. We examine baseline methods for 3D tracking with map-derived context. We also mine one thousand hours of fleet logs to find diverse, real-world object trajectories which constitute our motion forecasting benchmark. We examine baseline forecasting methods and see that map data significantly improves accuracy. We will maintain a public leaderboard for 3D object tracking and motion forecasting. The sensor data, map data, annotations, and code which make up Argoverse are available at [Argoverse.org](https://argoverse.org).

Acknowledgements. We thank our Argo AI colleagues – Ben Ballard, Brett Browning, Alex Bury, Dave Chekan, Kunal Desai, Patrick Gray, Larry Jackson, Etienne Jacques, Gang Pan, Kevin Player, Peter Rander, Bryan Salesky, Philip Tsai, Ian Volkwein, Ersin Yumer and many more – for their invaluable assistance in supporting Argoverse.

Patsorn Sangkloy is supported by a Royal Thai Government Scholarship. James Hays receives research funding from Argo AI, which is developing products related to the research described in

this paper. In addition, the author serves as a Staff Scientist to Argo AI. The terms of this arrangement have been reviewed and approved by Georgia Tech in accordance with its conflict of interest policies.

References

- [1] Alexandre Alahi, Kratharth Goel, Vignesh Ramanathan, Alexandre Robicquet, Li Fei-Fei, and Silvio Savarese. Social lstm: Human trajectory prediction in crowded spaces. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [2] Mykhaylo Andriluka, Stefan Roth, and Bernt Schiele. People-tracking-by-detection and people-detection-by-tracking. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2008.
- [3] Andrew Bacha, Cheryl Bauman, Ruel Faruque, Michael Fleming, Chris Terwelp, Charles Reinholtz, Dennis Hong, Al Wicks, Thomas Alberi, David Anderson, Stephen Cacciola, Patrick Currier, Aaron Dalton, Jesse Farmer, Jesse Hurdus, Shawn Kimmel, Peter King, Andrew Taylor, David Van Covern, and Mike Webster. Odin: Team victortango’s entry in the darpa urban challenge. *J. Field Robot.*, 25(8):467–492, Aug. 2008.
- [4] Keni Bernardin and Rainer Stiefelhagen. Evaluating multiple object tracking performance: The clear mot metrics. *EURASIP J. Image and Video Processing*, 2008.
- [5] Holger Caesar, Varun Bankiti, Alex H. Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for autonomous driving. *arXiv preprint arXiv:1903.11027*, 2019.
- [6] Sergio Casas, Wenjie Luo, and Raquel Urtasun. Intentnet: Learning to predict intention from raw sensor data. In Aude Billard, Anca Dragan, Jan Peters, and Jun Morimoto, editors, *Proceedings of The 2nd Conference on Robot Learning*, volume 87 of *Proceedings of Machine Learning Research*, pages 947–956. PMLR, 29–31 Oct 2018.
- [7] Florian Chabot, Mohamed Chaouch, Jaonary Rabarisoa, Céline Teulière, and Thierry Chateau. Deep manta: A coarse-to-fine many-task network for joint 2d and 3d vehicle analysis from monocular image. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [8] Chenyi Chen, Ari Seff, Alain Kornhauser, and Jianxiong Xiao. Deepdriving: Learning affordance for direct perception in autonomous driving. In *The IEEE International Conference on Computer Vision (ICCV)*, 2015.
- [9] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [10] Nachiket Deo and Mohan M Trivedi. Convolutional social pooling for vehicle trajectory prediction. *arXiv preprint arXiv:1805.06771*, 2018.
- [11] Martin Ester, Hans peter Kriegel, Jörg Sander, and Xiaowei Xu. A density-based algorithm for discovering clusters in large spatial databases with noise. In *KDD Proceedings of the Second International Conference on Knowledge Discovery and Data Mining*, pages 226–231. AAAI Press, 1996.
- [12] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The KITTI dataset. *The International Journal of Robotics Research*, 32(11):1231–1237, 2013.
- [13] Alex Graves. Generating sequences with recurrent neural networks. *CoRR*, abs/1308.0850, 2013.
- [14] Agrim Gupta, Justin Johnson, Li Fei-Fei, Silvio Savarese, and Alexandre Alahi. Social gan: Socially acceptable trajectories with generative adversarial networks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [15] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask r-cnn. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 2980–2988, Oct 2017.
- [16] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask R-CNN. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2017.
- [17] Simon Hecker, Dengxin Dai, and Luc Van Gool. End-to-end learning of driving models with surround-view cameras and route planners. In *European Conference on Computer Vision (ECCV)*, 2018.
- [18] David Held, Devin Guillory, Brice Rebsamen, Sebastian Thrun, and Silvio Savarese. A probabilistic framework for real-time 3d segmentation using spatial, temporal, and semantic cues. In *Proceedings of Robotics: Science and Systems*, 2016.
- [19] David Held, Jesse Levinson, and Sebastian Thrun. Precision tracking with sparse 3d and dense color 2d data. In *ICRA*, 2013.
- [20] David Held, Jesse Levinson, Sebastian Thrun, and Silvio Savarese. Combining 3d shape, color, and motion for robust anytime tracking. In *Proceedings of Robotics: Science and Systems*, Berkeley, USA, July 2014.
- [21] M. Himmelsbach and H. j. Wußlsche. Lidar-based 3d object perception. In *Proceedings of 1st International Workshop on Cognition for Technical Systems*, 2008.
- [22] Xinyu Huang, Xinjing Cheng, Qichuan Geng, Binbin Cao, Dingfu Zhou, Peng Wang, Yuanqing Lin, and Ruigang Yang. The apolloscape dataset for autonomous driving. In *arXiv:1803.06184*, 2018.
- [23] Simon Julier, Jeffrey Uhlmann, and Hugh F Durrant-Whyte. A new method for the nonlinear transformation of means and covariances in filters and estimators. *IEEE Transactions on automatic control*, 45(3):477–482, 2000.
- [24] Simon J Julier, Jeffrey K Uhlmann, and Hugh F Durrant-Whyte. A new approach for filtering nonlinear systems. In *Proceedings of 1995 American Control Conference-ACC’95*, volume 3, pages 1628–1632. IEEE, 1995.
- [25] Abhijit Kundu, Yin Li, and James M. Rehg. 3d-rcnn: Instance-level 3d object reconstruction via render-and-compare. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [26] Alex H. Lang, Sourabh Vora, Holger Caesar, Lubing Zhou, Jiong Yang, and Oscar Beijbom. Pointpillars: Fast encoders for object detection from point clouds. In *CVPR*, 2018.

- [27] Namhoon Lee, Wongun Choi, Paul Vernaza, Christopher Bongsoo Choy, Philip H. S. Torr, and Manmohan Krishna Chandraker. DESIRE: distant future prediction in dynamic scenes with interacting agents. *CoRR*, abs/1704.04394, 2017.
- [28] John Leonard, Jonathan How, Seth Teller, Mitch Berger, Stefan Campbell, Gaston Fiore, Luke Fletcher, Emilio Frazzoli, Albert Huang, Sertac Karaman, Olivier Koch, Yoshiaki Kuwata, David Moore, Edwin Olson, Steve Peters, Justin Teo, Robert Truax, Matthew Walter, David Barrett, Alexander Epstein, Keoni Maheloni, Katy Moyer, Troy Jones, Ryan Buckley, Matthew Antone, Robert Galejs, Siddhartha Krishnamurthy, and Jonathan Williams. A perception-driven autonomous urban vehicle. *J. Field Robot.*, 25(10):727–774, Oct. 2008.
- [29] Jesse Levinson, Jake Askeland, Jan Becker, Jennifer Dolson, David Held, Sören Kammel, J. Zico Kolter, Dirk Langer, Oliver Pink, Vaughan R. Pratt, Michael Sokolsky, Ganymed Stanek, David Michael Stavens, Alex Teichman, Moritz Werling, and Sebastian Thrun. Towards fully autonomous driving: Systems and algorithms. In *IEEE Intelligent Vehicles Symposium (IV), 2011, Baden-Baden, Germany, June 5-9, 2011*, pages 163–168, 2011.
- [30] Justin Liang and Raquel Urtasun. End-to-end deep structured models for drawing crosswalks. In *The European Conference on Computer Vision (ECCV)*, September 2018.
- [31] Xingyu Liu, Charles R Qi, and Leonidas J Guibas. FlowNet3D: Learning scene flow in 3d point clouds. *arXiv preprint arXiv:1806.01411*, 2019.
- [32] Wenjie Luo, Bin Yang, and Raquel Urtasun. Fast and furious: Real time end-to-end 3d detection, tracking and motion forecasting with a single convolutional net. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [33] Suraj M S, Hugo Grimmer, Lukas Platinsky, and Peter Ondruska. Visual vehicle tracking through noise and occlusions using crowd-sourced maps. In *Intelligent Robots and Systems (IROS), 2018 IEEE international conference on*, pages 4531–4538. IEEE, 2018.
- [34] Yuexin Ma, Xinge Zhu, Sibozhang, Ruigang Yang, Wenping Wang, and Dinesh Manocha. Trafficpredict: Trajectory prediction for heterogeneous traffic-agents. In *Proceedings of the 33rd National Conference on Artificial Intelligence, AAAI’19*. AAAI Press, 2019.
- [35] Will Maddern, Geoffrey Pascoe, Chris Linegar, and Paul Newman. 1 year, 1000 km: The Oxford Robotcar dataset. *The International Journal of Robotics Research*, 36(1):3–15, 2017.
- [36] A. Milan, L. Leal-Taixé, I. Reid, S. Roth, and K. Schindler. MOT16: A benchmark for multi-object tracking. *arXiv:1603.00831 [cs]*, Mar. 2016. arXiv: 1603.00831.
- [37] Michael Montemerlo, Jan Becker, Suhrud Bhat, Hendrik Dahlkamp, Dmitri Dolgov, Scott Ettinger, Dirk Haehnel, Tim Hilden, Gabe Hoffmann, Burkhard Huhne, Doug Johnston, Stefan Klumpp, Dirk Langer, Anthony Levandowski, Jesse Levinson, Julien Marcil, David Orenstein, Johannes Paefgen, Isaac Penny, Anna Petrovskaya, Mike Pflueger, Ganymed Stanek, David Stavens, Antone Vogt, and Sebastian Thrun. Junior: The stanford entry in the urban challenge. *J. Field Robot.*, 25(9):569–597, Sept. 2008.
- [38] nuTonomy. nuscenescenes. <https://www.nuscenes.org/>, 2018.
- [39] Gaurav Pandey, James R McBride, and Ryan M Eustice. Ford campus vision and lidar data set. *Int. J. Rob. Res.*, 30(13):1543–1552, Nov. 2011.
- [40] SeongHyeon Park, Byeongdo Kim, Chang Mook Kang, Chung Choo Chung, and Jun Won Choi. Sequence-to-sequence prediction of vehicle trajectory via LSTM encoder-decoder architecture. In *Intelligent Vehicles Symposium*, pages 1672–1678. IEEE, 2018.
- [41] Abhishek Patil, Srikanth Malla, Haiming Gang, and Yi-Ting Chen. The h3d dataset for full-surround 3d multi-object detection and tracking in crowded urban scenes. In *International Conference on Robotics and Automation*, 2019.
- [42] Luis Patino, Tom Cane, Alain Vallee, and James Ferryman. Pets 2016: Dataset and challenge. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 1–8, 2016.
- [43] Akshay Rangesh, Kevan Yuen, Ravi Kumar Satzoda, Rakesh Nattoji Rajaram, Pujitha Gunaratne, and Mohan M. Trivedi. A multimodal, full-surround vehicular testbed for naturalistic studies and benchmarking: Design, calibration and deployment. *CoRR*, abs/1709.07502, 2017.
- [44] Stephan R. Richter, Zeeshan Hayder, and Vladlen Koltun. Playing for benchmarks. In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 2232–2241, 2017.
- [45] R.B. Rusu and S. Cousins. 3d is here: Point cloud library (pcl). In *Robotics and Automation (ICRA), 2011 IEEE International Conference on*, pages 1–4, may 2011.
- [46] John P. Snyder. Map projections: A working manual. u.s. geological survey professional paper. page 61, 1987.
- [47] Xibin Song, Peng Wang, Dingfu Zhou, Rui Zhu, Chenye Guan, Yuchao Dai, Hao Su, Hongdong Li, and Ruigang Yang. ApolloCar3d: A large 3d car instance understanding benchmark for autonomous driving. *CoRR*, abs/1811.12222, 2018.
- [48] Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 3104–3112. Curran Associates, Inc., 2014.
- [49] Christopher Urmson, Joshua Anhalt, J. Andrew (Drew) Bagnell, Christopher R. Baker, Robert E. Bittner, John M. Dolan, David Duggins, David Ferguson, Tugrul Galatali, Hartmut Geyer, Michele Gittleman, Sam Harbaugh, Martial Hebert, Thomas Howard, Alonzo Kelly, David Kohanbash, Maxim Likhachev, Nick Miller, Kevin Peterson, Raj Rajkumar, Paul Rybski, Bryan Salesky, Sebastian Scherer, Young-Woo Seo, Reid Simmons, Sanjiv Singh, Jarrod M. Snider, Anthony (Tony) Stentz, William (Red) L. Whittaker, and Jason Ziegler. Tartan racing: A multi-modal approach to the darpa urban challenge. Technical report, Carnegie Mellon University, Pittsburgh, PA, April 2007.

- [50] Shenlong Wang, Min Bai, Gellert Mattyus, Hang Chu, Wenjie Luo, Bin Yang, Justin Liang, Joel Chaverie, Sanja Fidler, and Raquel Urtasun. Torontocity: Seeing the world with a million eyes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [51] Weiyue Wang, Ronald Yu, Qiangui Huang, and Ulrich Neumann. Sgpn: Similarity group proposal network for 3d point cloud instance segmentation. In *CVPR*, 2018.
- [52] Bin Yang, Ming Liang, and Raquel Urtasun. Hdnet: Exploiting hd maps for 3d object detection. In Aude Billard, Anca Dragan, Jan Peters, and Jun Morimoto, editors, *Proceedings of The 2nd Conference on Robot Learning*, volume 87 of *Proceedings of Machine Learning Research*, pages 146–155. PMLR, 29–31 Oct 2018.
- [53] Yin Zhou and Oncel Tuzel. Voxelnet: End-to-end learning for point cloud based 3d object detection. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.

7. Supplementary Material

In this supplementary material, we present additional details about our map (Section 8), our trajectory mining (Section 9), and our 3D tracking algorithm (Section 10).

8. Supplemental Map Details

In this section, we describe details of our map coordinate system and the functions exposed by our map API, and we visualize several semantic attributes of our vector map. Our map covers 204 linear kilometers of lane centerlines in Miami and 86 linear kilometers in Pittsburgh. In terms of driveable area, our map covers 788,510 m² in Miami and 286,104 m² in Pittsburgh.

8.1. Coordinate System

The model of the world that we subscribe to within our map and dataset is a local tangent plane centered at a central point located within each city. This model has a flat earth assumption which is approximately correct at the scale of a city. Thus, we provide map object pose values in city coordinates. City coordinates can be converted to the UTM (Universal Transverse Mercator) coordinate system by simply adding the city’s origin in UTM coordinates to the object’s city coordinate pose. The UTM model divides the earth into 60 flattened, narrow zones, each of width 6 degrees of longitude. Each zone is segmented into 20 latitude bands.

We favor a city-level coordinate system because of its high degree of interpretability when compared with geocentric reference coordinate systems such as the 1984 World Geodetic System (WGS84). While WGS84 is widely used by the Global Positioning System, the model is difficult to interpret at a city-scale; because its coordinate origin is located at the Earth’s center of mass, travel across an entire city corresponds only to pose value changes in the hundredth decimal place. The conversion back and forth between UTM and WGS84 is well-known and is documented in detail in [46].

We provide ground-truth object pose data in the ego-vehicle frame, meaning a single SE(3) transform is required to bring points

into the city frame for alignment with the map:

$$p_{city} = ({}^{city}T_{egovehicle}) (p_{egovehicle})$$

Figure 9 shows examples of the centerlines which are the basis of our vector map. Centerline attributes include whether or not lane segments are in an intersection, and which lane segments constitute their predecessors and successors. Figure 10 shows examples of centerlines, driveable area, and ground height projected onto camera image.

8.2. Map API and Software Development Kit

The dataset’s rich maps are our most significant contribution and we aim to make it easy to develop computer vision tools that leverage the map data. Figure 12 describes several functions which we hope will make it easier for researchers to access the map. Our API is provided in Python. For example, our API can provide rasterized bird’s eye view (BEV) images of the map around the ego-vehicle, extending up to 100 m in all directions. It can also provide a dense 1 meter resolution grid of the ground surface, especially useful for ground classification when globally planar ground surface assumptions are violated (see Figure 13).

These dense, pixel-level map renderings, similar to visualizations of instance-level or semantic segmentation [9], have recently been demonstrated to improve 3d perception and are relatively easy to use as an input to a convolutional network [52, 6].

We provide our vector map data in a modified OpenStreetMap (OSM) format, i.e. consisting of “Nodes” (waypoints) composed into “Ways” (polylines) so that the community can take advantage of open source mapping tools built to handle OSM formats. The data we provide is richer than existing OSM data which does not contain per-lane or elevation information.

9. Supplemental Details on Mined Trajectories for Forecasting

In this section, we describe our approach for mining data for trajectory forecasting. The scenarios challenging for a forecasting task are rare but with a vector map, they are easy to identify. We focus on some specific behavioral scenarios from over 1006 driving hours. For every 5 second sequence, we assign an *interesting* score to every track in that sequence. A high *interesting* score can be attributed to one or more of the following cases wherein the track is: at an intersection with or without traffic control, on a right turn lane, on a left turn lane, changing lanes to a left or right neighbor, having high median velocity, having high variance in velocity and visible for a longer duration. We give more importance to changing lanes and left/right turns because these scenarios are very rare. If there are at least 2 sufficiently important tracks in the sequence, we save the sequence for forecasting experiments. Further, the track which has the maximum *interesting* score and is visible through out the sequence is tagged as the *Agent*. The forecasting task is then to predict the trajectory of this particular track, where all the other tracks in the sequence can be used for learning social context for the *Agent*. There is also a 2.5 second overlap between 2 consecutive sequences. This means the same track id can be available in 2 sequences, albeit with different trajectories.

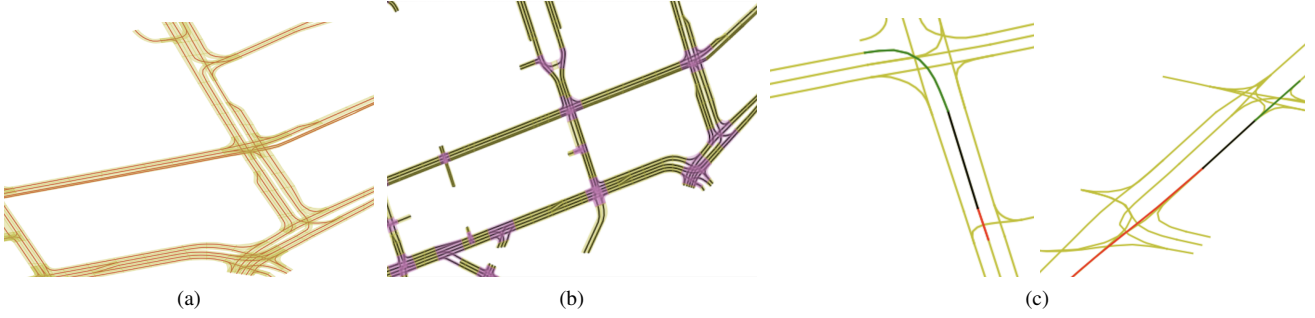


Figure 9: (a) Lane centerlines and hallucinated area are shown in red and yellow, respectively. We provide *lane* centerlines in our dataset because simple *road* centerline representations cannot handle the highly complicated nature of real world mapping, as shown above with divided roads. (b) We show lane segments within intersections in pink, and all other lane segments in yellow. Black shows lane centerlines. (c) Example of a specific lane centerline’s successors and predecessors. Red shows the predecessor, green shows the successor, and black indicates the centerline segment of interest.

10. Supplemental Tracking Details

In this section, we describe our tracking pipeline in greater detail.

10.1. Tracker Implementation Details

Because of space constraints we could not fit all details of our 3D tracking pipeline in the main paper. We do not claim any novelty for this ‘baseline’ tracker, but it works reasonably well, especially with map information to make the task easier (e.g. driveable area, ground height, and lane information). Some other recent attempts in 3D tracking and detection include H3D[41], which uses VoxelNet[53] for detection and an Unscented Kalman Filter [24, 23] for tracking. NuScenes dataset [38] uses PointPillar[26] for 3D bounding box detection. ApolloCar3D[47] implements 2D to 3D car pose estimation baselines based on 3D-RCNN[25] and DeepMANTA[7].

Our tracker tracks the position and velocity of surrounding vehicles from LiDAR data. The tracking pipeline has the following stages:

1. Segmentation and Detection. In order to segment a point cloud into distinct object instances, we exploit the complementary nature of our two sensor modalities. First, using Mask R-CNN [16], we obtain object masks in pixel space and discard any LiDAR returns whose image projection does not fall within a mask. We then geometrically cluster the remaining 3D LiDAR point cloud into separate objects according to density, using DBSCAN [11].

Others have proposed compensating for point cloud undersegmentation and oversegmentation scenarios by conditioning on the data association and then jointly track and perform probabilistic segmentation [18]. We can avoid many such segmentation failures with the high precision of our Mask R-CNN network². We also eliminate the need for an object’s full, pixel-colored 3D shape during tracking, as others have suggested [20, 19]. We prefer density-based clustering to connected components clustering in a 2D oc-

cupancy grid [21, 29] because the latter approach discards information along the z-axis, often rendering the method brittle.

To help focus our attention to areas that are important for a self driving car, we only consider points within the region of interest (ROI) defined by the driveable area map, and not on the ground.

While segmentation provides us a set of points belonging to an object, we need to determine if this is an object of interest that we want to track. Unlike in image space, objects in a 3D have consistent sizes. We apply heuristics that enforce the shape and volume of a typical car and thereby identify vehicle objects to be tracked. We estimate the center of an object by fitting a smallest enclosing circle over the segment points.

2. Association. We utilize the Hungarian algorithm to obtain globally optimal assignment of previous tracks to currently detected segments where the cost of assignment is based on spatial distance. Typically, tracks are simply assigned to their nearest neighbor in the next frame.

3. Tracking. We use ICP (Iterative Closest Point) from the Point Cloud Library [45] to estimate relative transformation between corresponding point segments for each track. Then we apply a Kalman Filter (KF) [19] with ICP results as the measurement and a static motion model (or constant velocity motion model, depending on the environment) to estimate vehicle poses for each tracked vehicle. We assign a fixed size bounding box for each tracked object. The KF state is comprised of both the 6 dof pose and velocity.

10.2. Tracking Evaluation Metrics

We use standard evaluation metrics commonly used for multiple object trackers (MOT) [36, 4]. The MOT metric relies on centroid distance as distance measure.

- **MOTA(Multi-Object Tracking Accuracy):**

$$MOTA = 100 * (1 - \frac{\sum_t FN_t + FP_t + ID_{sw}}{\sum_t GT}) \quad (1)$$

where FN_t , FP_t , ID_{sw} , GT denote number of false negative, false positives, number of ID switches, and ground truths. We report MOTA as percentages.

²We use a public implementation available at <https://github.com/facebookresearch/maskrcnn-benchmark>.



(a) Lane geometry and connectivity



(b) Driveable area



(c) Ground height

Figure 10: Examples of centerlines, driveable area, and ground height projected onto camera image

- **MOTP(Multi-Object Tracking Precision):**

$$MOTP = \frac{\sum_{i,t} D_t^i}{\sum_t C_t} \quad (2)$$

where C_t denotes the number of matches, and D_t^i denotes the distance of matches.

- **IDF1 (F1 score):**

$$IDF1 = 2 \frac{precision * recall}{precision + recall} \quad (3)$$

Where *recall* is the number of true positives over number of total ground truth labels. *precision* is the number of true positives over sum of true positives and false positives.

- **MT (Mostly Tracked):** the ratio of trajectories tracked more than 80% of its lifetime.

- **ML (Mostly Lost):** the ratio of trajectories tracked less than 20% of its lifetime.
- **FP (False Positive):** Total number of false positives
- **FN (False Negative):** Total number of false negatives
- **IDsw (ID Switch):** number of identified ID switches
- **Frag (Fragmentation):** Total number of switches from "tracked" to "not tracked"

10.3. True Positive Thresholding Discussion

Intersection-over-Union (IoU) is designed as a scale invariant metric, meaning that doubling the size and relative overlap of two boxes will not change its value. However, we counter that 3d tracking evaluation should not be performed in a strictly scale invariant manner. *Absolute* error matters, especially in 3d. In 2d tasks (e.g. object detection) we have only pixels which could be any real world size, whereas in 3d we have absolute lengths. When using IOU as a TP/FP threshold, with fixed intersection area, IoU for larger vehicles is penalized unfairly (see Figure 14). On the other hand, with a fixed distance between associated centroids, the IoU increases with larger vehicles. In the LiDAR domain, these problems are exaggerated because the sampling density can be quite low, especially for distant objects. In 2d object detection, we rarely try to find objects that are 3 pixels in size, but small, distant objects frequently have 3 LiDAR returns and thus accurate determination of their spatial extent is difficult.



Figure 11: **Ring Camera Examples.** Scenes captured in Miami, Florida, USA (top) and Pittsburgh, Pennsylvania, USA (bottom) with our ring camera. Each row consists of 7 camera views with overlapping fields of view. Camera order is *rear_left, side_left, front_left, front_center, front_right, side_right, rear_right*

Function name	Description
<code>remove_non_driveable_area_points</code>	Use rasterized driveable area ROI to decimate LiDAR point cloud to only ROI points.
<code>remove_ground_surface</code>	Remove all 3D points within 30 cm of the ground surface.
<code>get_ground_height_at_xy</code>	Get ground height at provided (x,y) coordinates.
<code>render_local_map_bev_cv2</code>	Render a Bird's Eye View (BEV) in OpenCV.
<code>render_local_map_bev_mpl</code>	Render a Bird's Eye View (BEV) in Matplotlib.
<code>get_nearest_centerline</code>	Retrieve nearest lane centerline polyline.
<code>get_lane_direction</code>	Retrieve most probable tangent vector $\in \mathbb{R}^2$ to lane centerline.
<code>get_semantic_label_of_lane</code>	Provide boolean values regarding the lane segment, including <i>is_intersection</i> , <i>turn_direction</i> , and <i>has_traffic_control</i> .
<code>get_lane_ids_in_xy_bbox</code>	Get all lane IDs within a Manhattan distance search radius in the <i>xy</i> plane.
<code>get_lane_segment_predecessor_ids</code>	Retrieve all lane IDs with an incoming edge into the query lane segment in the semantic graph.
<code>get_lane_segment_successor_ids</code>	Retrieve all lane IDs with an outgoing edge from the query lane segment.
<code>get_lane_segment_adjacent_ids</code>	Retrieve all lane segment IDs of that serve as left/right neighbors to the query lane segment.
<code>get_lane_segment_centerline</code>	Retrieve polyline coordinates of query lane segment ID.
<code>get_lane_segment_polygon</code>	Hallucinate a lane polygon based around a centerline using avg. lane width.
<code>get_lane_segments_containing_xy</code>	Use a "point-in-polygon" test to find lane IDs whose hallucinated lane polygons contain this (x, y) query point.

Figure 12: Example Python functions in the Argoverse map API.



Figure 13: A scene with non-planar ground surface. The colored LiDAR returns have been classified as ground based on the map. Points outside the drivable area are also discarded. This simple distance threshold against a map works well, even on the road to the left which goes steeply uphill.

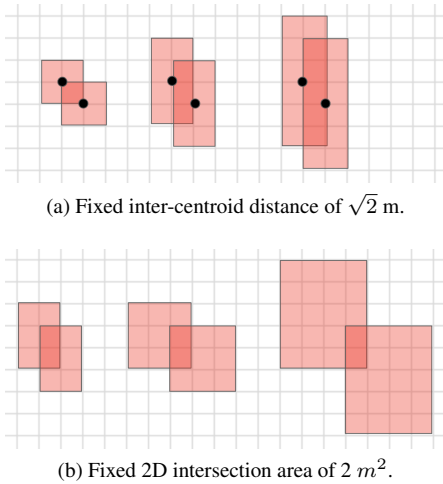


Figure 14: We compare thresholding true positives (TP) and false positives (FP) of matched cuboid shapes using inter-centroid distance (above) versus using 2D/3D IoU (below). Above: fixed inter-centroid distance, from left to right: IoU values of 0.143, 0.231, 0.263. Below: fixed intersection area, from left to right, IoU values of 0.2, 0.125, 0.053.

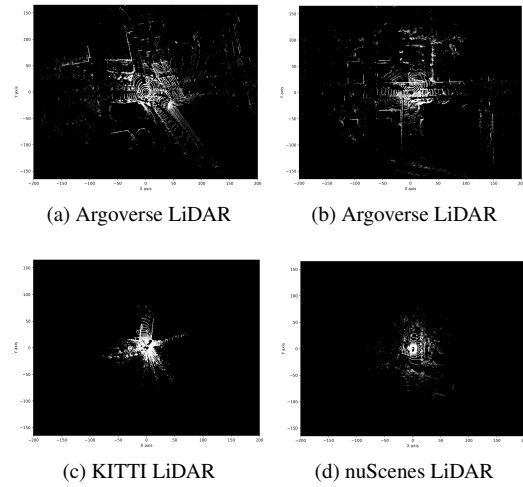


Figure 15: Above: Sample LiDAR sweeps in the ego-vehicle frame, with marked x and y axes, with $x \in [-200, 200]$ and $y \in [-160, 160]$ for all plots. The Argoverse LiDAR has twice the range of the sensors used to collect the KITTI or nuScenes datasets, allowing us to observe more objects in each scene.