

# Argoverse: 3D Tracking and Forecasting with Rich Maps

Ming-Fang Chang<sup>\*1,2</sup>, John Lambert<sup>\*3</sup>, Patsorn Sangkloy<sup>\*1,3</sup>, Jagjeet Singh<sup>\*1</sup>, Sławomir Bąk, Andrew Hartnett<sup>1</sup>, De Wang<sup>1</sup>, Peter Carr<sup>1</sup>, Simon Lucey<sup>1,2</sup>, Deva Ramanan<sup>1,2</sup>, and James Hays<sup>1,3</sup>

<sup>1</sup>Argo AI, <sup>2</sup>Carnegie Mellon University, <sup>3</sup>Georgia Institute of Technology

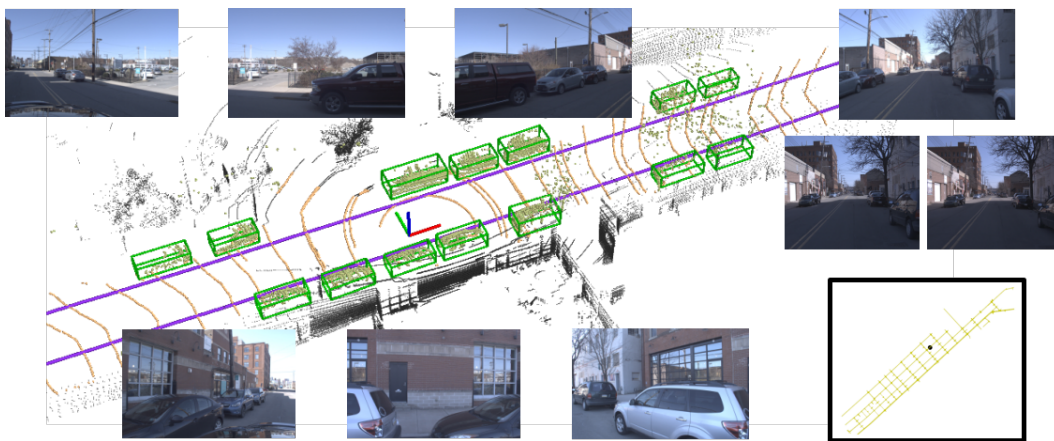


Figure 1: We introduce datasets for 3D tracking and motion forecasting with *rich maps* for autonomous driving. Our 3D tracking dataset contains sequences of LiDAR measurements, 360° RGB video, front-facing stereo (middle-right), and 6-dof localization. All sequences are aligned with maps containing 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 – two datasets designed to support autonomous vehicle machine learning tasks such as 3D tracking and motion forecasting. Argoverse was collected by a fleet of autonomous vehicles in Pittsburgh and Miami. The Argoverse 3D Tracking dataset includes 360° images from 7 cameras with overlapping fields of view, 3D point clouds from long range LiDAR, 6-DOF pose, and 3D track annotations. Notably, it is the only modern AV dataset that provides forward-facing stereo imagery. The Argoverse Motion Forecasting dataset includes more than 300,000 5 second tracked scenarios with a particular vehicle identified for trajectory forecasting. Argoverse is the first autonomous vehicle dataset to include “HD maps” with 290 km of mapped lanes with geometric and semantic metadata. All data is released under a Creative Commons license at

[www.argoverse.org](http://www.argoverse.org). In our baseline experiments, we illustrate how detailed map information such as lane direction, driveable area, and ground height improves the accuracy of 3D object tracking and motion forecasting. Our tracking and forecasting experiments represent only an initial exploration of the use 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 of KITTI [14], which opened and connected a plethora of new research directions. However, publicly available datasets for autonomous driving rarely include *map* data, even though detailed maps are

\*Equal contribution

critical to the development of 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, which lane a tracked object is in, how far it is to the next intersection, etc. But since publicly available datasets do not contain richly-mapped attributes, *how* to represent and utilize such features is an open research question. Argoverse is the first large-scale autonomous driving dataset with such detailed maps. We investigate 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 3D tracking dataset with synchronized data from LiDAR, 360° and stereo cameras sampled across two cities in varied conditions. Unlike other recent datasets, our 360° is captured at 30fps.
- We provide ground truth 3D track annotations across 15 object classes, with five times as many tracked objects as the KITTI [14] tracking benchmark.
- We create a large-scale forecasting dataset consisting of trajectory data for interesting scenarios such as turns at intersections, high traffic clutter, and lane changes.
- We release map data and an API which can be used to develop map-based perception and forecasting algorithms. We are the first self-driving vehicle dataset with a semantic vector map of road infrastructure and traffic rules. The inclusion of “HD” map information also means our dataset is the first large-scale benchmark for automatic map creation, often known as *map automation*.
- We are the first to examine the influence of HD map context for 3D tracking and motion forecasting. In the case of 3D tracking, we measure the influence of map-based ground point removal and orientation snapping to lanes. In the case of motion forecasting, we experiment with the creation of diverse predictions from the lane graph and the pruning of predictions by the driveable area map. In both cases, we see higher accuracy with the use of a map.

## 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. The prohibitive cost of annotating and constructing such maps has spurred interest in the growing field of map automation [35, 25, 4]. Prior to Argoverse’s release, no public dataset

included 3D vector map information, thus preventing the development of common benchmark for map automation. TorontoCity [58] also focuses on map construction tasks but without 3D annotation for dynamic objects. The nuScenes dataset [6] originally contained 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. Subsequent to Argoverse release, nuScenes has released labels for 2D semantic map regions, without a lane or graph structure. Like nuScenes, we include maps of driveable area, but also include ground height and a “vector map” of lane centerlines and their connectivity.

### Autonomous Driving Datasets with 3D Track Annotations.

Many existing datasets for object tracking focus on pedestrian tracking from image/video sequences [16, 48, 43, 2]. Several datasets provide raw data from self-driving vehicle sensors, but without any object annotations [42, 45, 49]. The ApolloCar3D dataset [55] is oriented towards 3D semantic object keypoint detection instead of tracking. KITTI [14] and H3D [47] offer 3D bounding boxes and track annotations but do not provide a map. The camera field of view is frontal, rather than 360°. VIPER [52] provides data from a simulated world with 3D track annotations. nuScenes [6] currently provides 360° data and a benchmark for 3D object detection, with tracking annotation also available. The *Argoverse 3D 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 Trajectory Data.

ApolloScape [26] also uses sensor-equipped vehicles to observe driving trajectories in the wild and presents a forecasting benchmark [41] from a subset of the ApolloScape 3D tracking annotations. This dataset consists of 155 minutes of observations compared to 320 hours of observations in the *Argoverse Forecasting* dataset. IntentNet [7] mines roof-mounted LiDAR data for 54 million object trajectories, but the data is not publicly available.

### 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 [60] 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* (RNDF) [44, 57, 3]. Hecker *et al.* [20] show that end-to-end route planning can be improved by processing rasterized maps from OpenStreetMap and TomTom. Liang *et al.* [36] demonstrate that using road centerlines and intersection polygons from OpenStreetMap can help infer crosswalk location and direction. Yang *et al.* [60] show that incorporating ground height and bird’s eye view (BEV) road

segmentation with LiDAR point information as a model input can improve 3D object detection. Liang *et al.* [37] show how 3D object detection accuracy can be improved by using mapping (ground height estimation) as an additional task in multi-task learning. Suraj *et al.* [40] 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 [13, 33] or connected components of an occupancy grid [34, 24], and then associated based on some distance function using the Hungarian algorithm. Held *et al.* utilize probabilistic approaches to point cloud segmentation and tracking [21, 23, 22]. 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 [59, 38]. 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 [17] attempts to model the multimodal nature of the predictions. However, both have only been tested on pedestrian trajectories, with no use of static context (e.g. a map). Deo *et al.* [11] 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.* [9] 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 [7] extends the joint detection and prediction approach of Luo *et al.* [39] by discretizing the prediction space and attempting to predict one of eight common driving maneuvers. DESIRE [32] 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. Another challenging problem in the trajectory forecasting domain is to predict diverse trajectories which can address multimodal nature of the problem. R2P2 [50] address the diversity-precision trade-off of generative forecasting models and formulate a symmetric cross-entropy training objective to address it. It is then followed by PRECOG [51] wherein they present the first generative multi-agent forecasting method to condition on agent intent. They achieve state-of-the-art results for forecasting methods in real (nuScenes [6]) and simulated

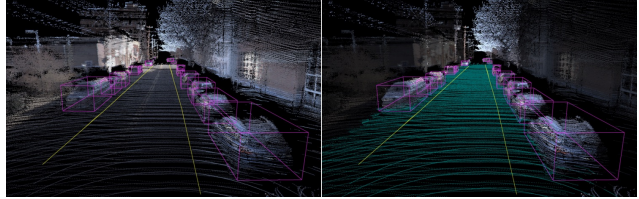


Figure 2: **3D visualization of an Argoverse scene.** Left: we accumulate LiDAR points and project them to a virtual image plane. Right: using our map, LiDAR points beyond driveable area are dimmed and points near the ground are highlighted in cyan. Cuboid object annotations and road centerlines are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

(CARLA [12]) datasets.

### 3. The Argoverse Dataset

Our sensor data, maps, and annotations are the primary contribution of this work. We also provide an API which connects the map data with sensor information *e.g.* ground point removal, nearest centerline queries, and lane graph connectivity; see the Appendix for more details. The data is available at [www.argoverse.org](http://www.argoverse.org) under a Creative Commons license. The API, tutorials, and code for baseline algorithms are available at [github.com/argoai/argoverse-api](https://github.com/argoai/argoverse-api) under an MIT license. The statistics and experiments in this document are based on Argoverse v1.1 released in October 2019.

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

**Sensors.** Our vehicles are equipped with two roof-mounted, rotating 32 beam LiDAR sensors. Each LiDAR has a 40° vertical field of view, with 30° overlapping field of view and 50° total field of view with both LiDAR. LiDAR range is up to 200 meters, roughly twice the range 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 [6] dataset (ours  $\sim 107,000$  points vs. nuScenes  $\sim 35,000$  points). The two LiDAR sensors rotate at 10 Hz and are out of phase, *i.e.* rotating in the same direction and speed but with an offset to avoid interference. 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  $\times$  1200) recording at 30 Hz with overlapping fields of view, providing 360°

| DATASET NAME                                    | MAP TYPE      | EXTENT OF ANNOTATED LANES    | DRIVEABLE AREA COVERAGE  | CAMERA FRAME RATE | 360° CAMERAS | INCLUDES STEREO | # TRACKED OBJECTS /SCENE | # SCENES |
|---|---------------|------------------------------|--------------------------|-------------------|--------------|-----------------|--------------------------|----------|
| KITTI [14]                                      | None          | 0 km                         | 0 m <sup>2</sup>         | 10 Hz             | no           | ✓               | 43.67 (train)            | 50       |
| Oxford RobotCar [42]                            | None          | 0 km                         | 0 m <sup>2</sup>         | 11/16Hz           | no           | no              | 0                        | 100+     |
| H3D [47]  | None          | 0 km                         | 0 m <sup>2</sup>         | 30 Hz             | no           | no              | 86.02 (train+val+test)   | 160      |
| Lyft Dataset [29] <sup>1</sup>                  | Raster        | 0 km                         | 48,690 m <sup>2</sup>    | 10 Hz             | ✓            | no              | 102.34 (train)           | 180+     |
| nuScenes v1.0 [6]                               | Vector+Raster | 133 km                       | 1,115,844 m <sup>2</sup> | 12 Hz             | ✓            | no              | 75.75 (train+val)        | 1000     |
| ApolloScape Tracking [41]                       | None          | 0 km                         | 0 m <sup>2</sup>         | n/a               | no           | no              | 206.16 (train)           | 103      |
| Waymo Open Dataset                              | None          | 0 km                         | 0 m <sup>2</sup>         | 10 Hz             | ✓            | no              | 113.68 (train+val)       | 1000     |
| Argoverse 3D Tracking v1.1 (human annotated)    | Vector+Raster | 204 km (MIA)<br>+86 km (PIT) | 1,192,073 m <sup>2</sup> | 30 Hz             | ✓            | ✓               | 97.81 (train+val+test)   | 113      |
| ApolloScape Forecasting [41]                    | None          | 0 km                         | 0 m <sup>2</sup>         | n/a               | no           | no              | 50.06 (train)            | 103      |
| Argoverse Forecasting v1.1 (mined trajectories) | Vector+Raster | 204 km (MIA)<br>+86 km (PIT) | 1,192,073 m <sup>2</sup> | -                 | no           | no              | 50.03 (train+val+test)   | 324,557  |

Table 1: **Public self-driving datasets.** We compare recent, publicly available self-driving datasets with 3D object annotations for tracking (top) and trajectories for forecasting (bottom). Coverage area for nuScenes is based on its *road and sidewalk* raster map. Argoverse coverage area is based on our *driveable area* raster map. Statistics updated September 2019.

coverage. In addition, there are 2 front-facing stereo cameras (2056 × 2464 with a 0.2986 m baseline) 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 Appendix. 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 [49], we place the origin of the ego-vehicle coordinate system at the center of the rear axle. All LiDAR data is provided in the ego-vehicle coordinate system, rather than in the respective LiDAR sensor coordinate frames. 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

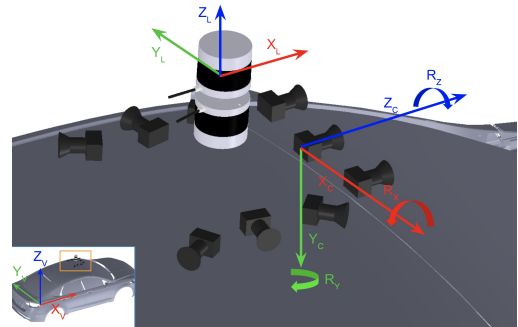


Figure 3: **Car sensor schematic.** Three reference coordinate systems are displayed: (1) the *vehicle frame*, with  $X_v$  forward,  $Y_v$  left, and  $Z_v$  up, (2) the *camera frame*, with  $X_c$  across image plane,  $Y_c$  down image plane, and  $Z_c$  along optical axis, (3) the *LiDAR frame*, 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.

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 are 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 there can be multiple (for merges and splits, respectively).

Centerlines are provided as “polylines”, *i.e.* an ordered sequence of straight segments. Each straight segment is defined by 2 vertices:  $(x_i, y_i, z_i)$  start and  $(x_{i+1}, y_{i+1}, z_{i+1})$  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  $3.84 \pm 0.89$  m wide. In Pittsburgh, the average width is  $3.97 \pm 1.04$  m. Other types of lane segments that would not be suitable for self-driving, *e.g.* bike lanes, can be as narrow as 0.97 m in Miami and as narrow as 1.06 m 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. We annotate 3D objects with track labels if they are within 5 meters of the driveable area (Section 3.2). 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 grid 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 shows a cross section of a scene with uneven ground height.

### 3.2. 3D Track Annotations

The *Argoverse Tracking Dataset* contains 113 vehicle log segments with human-annotated 3D tracks. These 113 segments vary in length from 15 to 30 seconds and collectively contain 11,052 tracked objects. We compared these with other datasets in Table 1. For each log segment, we annotated all objects of interest (both dynamic and static) with bounding cuboids which follow the 3D LiDAR returns associated with each object over time. We only annotated objects within 5 m 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. The same object ID is used for the same object, even if temporarily occluded. Each object is labeled with one of 15 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 shows 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 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 mine for vehicles that are either (1) at intersections, (2) taking left or right turns, (3) changing to adjacent lanes, or (4) in dense traffic. In total, we collect 324,557 five second sequences and use them in the forecasting benchmark. Figure 6 shows the geographic distribution of these sequences. Each sequence contains the 2D, bird’s eye view centroid of each tracked object sampled at 10 Hz. The “focal” object in each sequence is always a vehicle, but the other tracked objects can be vehicles, pedestrians, or bicycles. Their trajectories are available as context for “social” forecasting models. The 324,557 sequences are split into 205,942 train, 39,472 validation, and 78,143 test sequences. Each sequence has one challenging trajectory which is the focus of our forecasting benchmark. The train, validation, 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. 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 investigate how various baseline tracking methods perform on the Argoverse 3D tracking benchmark. Our baseline methods utilize a hybrid approach with LiDAR and ring camera images and operate directly in 3D. In addition to measuring the baseline difficulty of our benchmark, we measure how simple map-based heuristics can influence tracking accuracy. For these baselines, we track and evaluate *vehicles* only.

Given a sequence of  $F$  frames, where each frame contains a set of ring camera images and 3D points from Li-



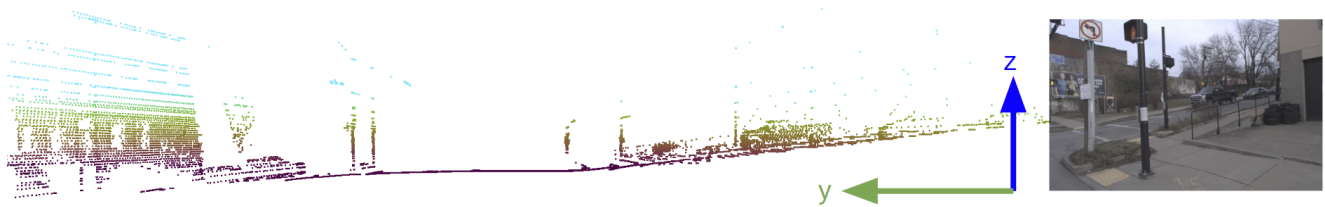


Figure 4: **Uneven ground scene in Argoverse dataset.** Some Argoverse scenes contain uneven ground, which is challenging to remove with simple heuristics (e.g. assuming that ground is planar). Above, we show a LiDAR slice with a slope on the right side and corresponding front-right camera image.

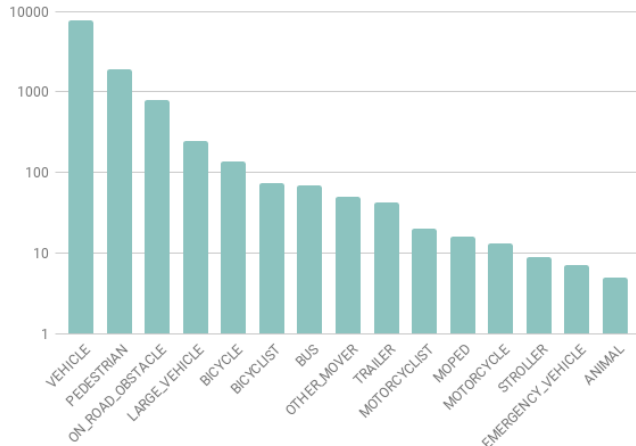


Figure 5: **Distribution of object classes.** This plot shows, in log scale, the number of 3D object tracks annotated for each class in the 113 log segments in the Argoverse 3D Tracking dataset.

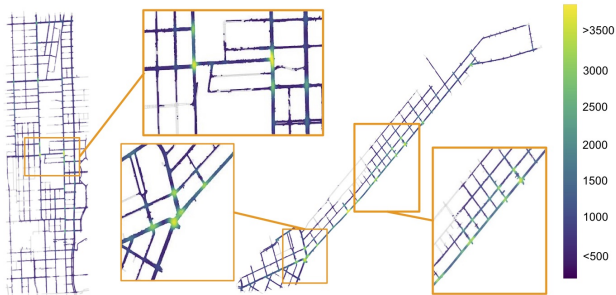


Figure 6: **Distribution of mined trajectories.** The colors indicate the number of mined trajectories across the maps of Miami (left) and Pittsburgh (right). The heuristics to find interesting vehicle behavior lead to higher concentrations in intersections and on busy roads such as Liberty and Penn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

DAR  $\{P_i \mid i = 1, \dots, N\}$ , where  $P_i \in \mathbb{R}^3$  of  $x, y, z$  coordinates, we want to determine a set of track hypotheses  $\{T_j \mid j = 1, \dots, n\}$  where  $n$  is the number of unique objects

in the whole sequence, and  $T_j$  contains the set of object center locations and orientation. 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.

**Baseline Tracker.** Our baseline tracking pipeline clusters LiDAR returns in driveable region (labeled by the map) to detect potential objects, uses Mask R-CNN [18] to prune non-vehicle LiDAR returns, associates clusters over time using nearest neighbor and the Hungarian algorithm, estimates transformations between clusters with iterative closest point (ICP), and estimates vehicle pose with a classical Kalman Filter using constant velocity motion model. The same predefined bounding box size is used for all vehicles.

When no match can be found by Hungarian method for an object, the object pose is maintained using only motion model up to 5 frames before being removed or associated to a new cluster. This enables our tracker to maintain same object ID even if the object is occluded for a short period of time and reappears. If a cluster is not associated with current tracked objects, we initialize a new object ID for it.

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 driveable area covers any region where it is possible for the vehicle to drive (see Section 3.1). This constraint reduces the opportunities for false positives.

**Ground height.** We use map information to remove LiDAR returns on the ground. 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 leverage standard evaluation metrics commonly used for multiple object tracking (MOT) [43, 5]. The MOT metric relies on a distance/similarity metric between ground truth and predicted objects to determine an optimal assignment. Once an assignment is made, we use three distance metrics for MOTP: MOTP-D (centroid distance), MOTP-O (orientation error), and MOTP-I (Intersection-over-Union error). MOTP-D is computed by the 3D bounding box centroid distance between associated tracker output and ground truth, which is also used in MOTA as detection association range. Our threshold for “missed” tracks is 2 meters, which is half of the average family car length in the US. MOTP-O is the smallest angle difference about the z (vertical) axis such that the front/back object orientation is ignored, and MOTP-I is the amodal shape estimation error, computed by the  $1 - IoU$  of 3D bounding box after aligning orientation and centroid as in nuScenes [6]. For all three MOTP scores, lower scores indicate higher accuracy.

In our experiments, we run our tracker over the 24 logs in the *Argoverse 3D Tracking* 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 baseline tracker performs well at short range where the LiDAR sampling density is higher, but struggles for objects beyond 50 m.

We compare our baseline tracker with three ablations that include: 1) using map-based ground removal and lane direction from the map; 2) using naive plane-fitting ground removal and lane direction from the map; 3) using map-based ground removal and no lane direction from the map. The results in Table 3 show that map-based ground removal leads to better 3D IoU score and slightly better detection performance (higher MOTA) than a plane-fitting approach at longer ranges, but slightly worse orientation. On the other hand, lane direction from the map significantly improves orientation performance, as shown in Figure 7.

We have employed relatively simple baselines to track objects in 3D. We believe that our data enables new approaches to map-based and multimodal tracking research.

## 5. Motion Forecasting

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

**1. Preprocessing:** As described in Section 3.3, we first mine for “interesting” sequences where a “focal” vehicle is observed for 5 seconds. As context, we have the centroids of all other tracked objects (including the AV itself) which are collapsed into one “other” class.

### Forecasting Coordinate System and Normalization.

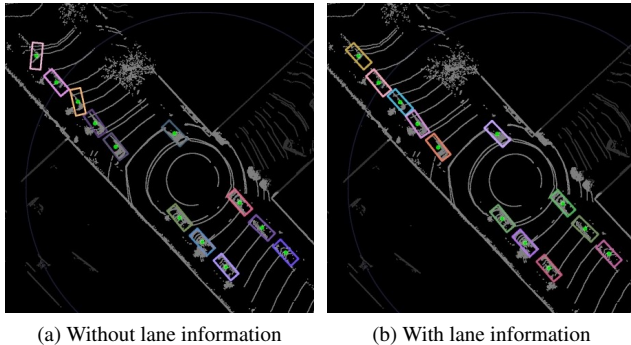


Figure 7: **Tracking with orientation snapping.** Using lane direction information helps to determine the vehicle orientation for detection and tracking.

The coordinate system we used for trajectory forecasting is a top-down, bird’s 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 C of the Appendix). (2) For models using lane centerlines as a reference path, we defined a *2D curvilinear coordinate system* with axes tangential and perpendicular to the lane centerline. (3) For models without the reference path (without a map), we normalize trajectories 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 normalization makes sure that  $y_i^{T_{obs}} = 0$ , 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 to capture social or spatial context. For social context, we use the 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 use the map as a prior by computing features 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 the 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 Constant Velocity, Nearest Neighbor, and LSTM Encoder-Decoder based [46, 15, 56] 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  $V_i = (X_i, Y_i)$  where*

| RANGE (m) | MOTA        | MOTP-D      | MOTP-O      | MOTP-I      | IDF1        | MT(%)       | ML(%)       | #FP         | #FN          | IDSW       | #FRAG      |
|-----------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|--------------|------------|------------|
| 30        | <b>65.5</b> | <b>0.71</b> | 15.3        | <b>0.25</b> | <b>0.71</b> | <b>0.67</b> | <b>0.18</b> | <b>5739</b> | <b>10098</b> | <b>356</b> | <b>380</b> |
| 50        | 50.0        | 0.81        | 13.5        | 0.26        | 0.59        | 0.30        | 0.31        | 8191        | 30468        | 607        | 691        |
| 100       | 34.2        | 0.82        | <b>13.3</b> | <b>0.25</b> | 0.46        | 0.13        | 0.51        | 9225        | 66234        | 679        | 773        |

Table 2: **Tracking accuracy at different ranges using map for ground removal and orientation initialization.** From top to bottom, accuracy for vehicles within 30 m, 50 m, and 100 m.

| RANGE (m) | USE MAP LANE | GROUND REMOVAL | MOTA        | MOTP-D (m)  | MOTP-O (degree) | MOTP-I      |
|-----------|--------------|----------------|-------------|-------------|-----------------|-------------|
| 30        | Y            | map            | 65.5        | <b>0.71</b> | 15.3            | <b>0.25</b> |
|           | Y            | plane-fitting  | <b>65.8</b> | 0.72        | <b>13.7</b>     | 0.29        |
|           | N            | map            | 65.4        | <b>0.71</b> | 25.3            | <b>0.25</b> |
| 50        | Y            | map            | <b>50.0</b> | <b>0.81</b> | 13.5            | <b>0.26</b> |
|           | Y            | plane-fitting  | 49.3        | <b>0.81</b> | <b>12.5</b>     | 0.29        |
|           | N            | map            | 49.8        | <b>0.81</b> | 27.7            | <b>0.26</b> |
| 100       | Y            | map            | <b>34.2</b> | <b>0.82</b> | 13.3            | <b>0.25</b> |
|           | Y            | plane-fitting  | 33.6        | <b>0.82</b> | <b>12.5</b>     | 0.28        |
|           | N            | map            | 34.1        | <b>0.82</b> | 27.7            | <b>0.25</b> |

Table 3: **Tracking accuracy comparison at different ranges while using different map attributes.** From top to bottom, accuracy for vehicles within 30 m, 50 m, and 100 m.

$X_i = (x_i^t, y_i^t)$  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 s is sufficient to capture the salient part of a trajectory, e.g. crossing an intersection. In this paper, we define the motion forecasting task as observing 20 past frames (2 s) and then predicting 30 frames (3 s) into the future. Each forecasting task can leverage the trajectories of other objects in the same sequence to capture the social context and map information for spatial context.

## 5.2. Evaluation of Multiple Forecasts

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 similar to DESIRE [32], Social GAN [17], R2P2 [50] and [51] wherein they encourage algorithms to output multiple predictions. Among the variety of metrics evaluated in [50] was the *minMSD* over  $K$  number of samples metric, where  $K = 12$ . A similar metric is used in [32] where they allow  $K$  to be up to 50. We follow the same approach and use minimum Average Displacement Error (minADE) and minimum Final Displacement Error (minFDE) over  $K$  predictions as our metrics, where  $K = 1, 3, 6, 9$ . Note that minADE refers to ADE of the trajectory which has minimum FDE, and not minimum ADE, since we want to evaluate the single best forecast.

That said, minADE error might not be a sufficient metric. As noted in [50] and [51], metrics like minMSD or minFDE can only evaluate how good is the best trajectory, but not how good are all the trajectories. A model having 5 good trajectories will have the same error as the model having 1 good and 4 bad trajectories. Further, given the multimodal nature of the problem, it might not be fair to evaluate against a single ground truth. In an attempt to evaluate based on the quality of predictions, we propose another metric: Drivable Area Compliance (DAC). If a model produces  $n$  possible future trajectories and  $m$  of those exit the drivable area at some point, the DAC for that model would be  $(n - m)/n$ . Hence, higher DAC means better quality of forecasted trajectories. Finally, we also use Miss Rate (MR) [61] with a threshold of 1.0 meter. It is again a metric derived from the distribution of final displacement errors. If there are  $n$  samples and  $m$  of them had the last coordinate of their best trajectory more than 2.0 m away from ground truth, then miss rate is  $m/n$ . 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. We use centerlines as a form of hypothetical reference paths for the future. Our heuristics generate  $K = 10$  centerlines. Our map gives us an easy way to produce a compact yet diverse set of forecasts. Nearest Neighbor baselines can further predict variable number of outputs by considering different number of neighbors.



### 5.3. Results

In this section, we evaluate the effect of multimodal predictions, social context, and spatial context (from the vector map) to improve motion forecasting over horizons of 3 seconds into the future. We evaluated 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*: Nearest Neighbor regression where trajectories are queried by  $(x_i^t, y_i^t)$  for  $t = \{1, \dots, T_{obs}\}$ . To make  $K$  predictions, we performed a lookup for  $K$  Nearest Neighbors.
- *LSTM*: 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  $t = \{T_{obs}+1, \dots, T_{pred}\}$ . This is limited to one prediction because we used a deterministic model.
- *LSTM+social*: Similar to *LSTM* but with input as  $(x_i^t, y_i^t, s_i^t)$ , where  $s_i^t$  denotes social features.
- *NN+map(prune)*: This baseline builds on *NN* and prunes the number of predicted trajectories based on how often they exit the drivable area. Accordingly, this method prefers predictions which are qualitatively good, and not just Nearest Neighbors.
- *NN+map(prior) m-G,n-C*: Nearest Neighbor regression where trajectories are queried by  $(a_i^t, o_i^t)$  for  $t = \{1, \dots, T_{obs}\}$ . m-G, n-C refers to  $m$  guesses (m-G) allowed along each of  $n$  different centerlines (n-C). Here,  $m > 1$ , except when  $K = 1$ .
- *NN+map(prior) 1-G,n-C*: This is similar to the previous baseline. The only difference is that the model can make only 1 prediction along each centerline.
- *LSTM+map(prior) 1-G,n-C*: Similar to *LSTM* 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 the centerlines. Distances  $(a_i^t, o_i^t)$  are then mapped to  $(x_i^t, y_i^t)$  for evaluation. Further, we make only one prediction along each centerline because we used a deterministic model.

The results of these baselines are reported in Table 4. When only 1 prediction is allowed, NN based baselines suffer from inaccurate neighbors and have poor minADE and minFDE. On the other hand, LSTM based baselines are able to at least learn the trajectory behaviors and have better results. *LSTM* baselines with no map are able to obtain the best minADE and minFDE for  $K = 1$ . Also, baselines which use map as prior have a much higher DAC. Now, as  $K$  increases, *NN* benefits from the map prior and consistently produces better predictions. When map is used for pruning, it further improves the selected trajectories and provides the best minADE and minFDE. *LSTM+map(prior) 1-G,n-C* outperforms *NN+map(prior) 1-G,n-C* highlighting the fact that LSTM does a better job generalizing to curvilinear coordinates. Further, using the

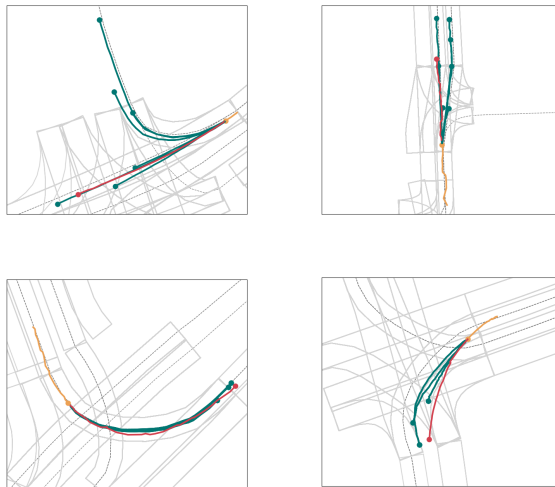


Figure 8: **Qualitative results from *NN+map(prior) m-G,n-C* motion forecasting baseline.** The orange trajectory represents the observed 2 s. Red represents ground truth for the next 3 seconds and green represents the multiple forecasted trajectories for those 3 s. **Top left:** The car starts to accelerate from a stop line and the model is able to predict 2 different modes (right turn and go straight) along with different velocity profiles along those modes. **Top right:** The model is able to predict 2 different scenarios – lane change and staying in the same lane. **Bottom left:** The model is able to cross a complex intersection and take a wide left turn without violating any lane rules because it is able to use the vector map to generate a reference path. **Bottom right:** The predictions account for different ways in which a left turn can be taken, in terms of velocity profile and turning radius.

map as a prior always provides better DAC, proving that our map helps in forecasting trajectories that follow basic map rules like staying in the driveable area. Another interesting comparison is between *NN+map(prior) 1-G,n-C* and *NN+map(prior) m-G,n-C*. The former comes up with many reference paths (centerlines) and makes one prediction along each of those paths. The latter comes up with fewer reference paths but produces multiple predictions along each of those paths. The latter outperforms the former in all 3 metrics, showing the importance of predicting trajectories which follow different velocity profiles along the same reference paths. Figure 9 reports the results of an ablation study for different values of  $m$  and  $n$ . Finally, when having access to HD vector maps and being able to make multiple predictions ( $K = 6$ ), even a shallow model like *NN+map(prior) m-G,n-C* is able to outperform a deterministic deep model *LSTM+social* ( $K = 1$ ) which has access to social context.

| BASELINE                | K=1         |             |             |             | K=3         |             |             |      | K=6         |             |             |             |
|-------------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|------|-------------|-------------|-------------|-------------|
|                         | minADE ↓    | minFDE ↓    | DAC ↑       | MR ↓        | minADE ↓    | minFDE ↓    | DAC ↑       | MR ↓ | minADE ↓    | minFDE ↓    | DAC ↑       | MR ↓        |
| Constant Velocity       | 3.53        | 7.89        | 0.88        | 0.83        | -           | -           | -           | -    | -           | -           | -           | -           |
| NN                      | 3.45        | 7.88        | 0.87        | 0.87        | 2.16        | 4.53        | 0.87        | 0.70 | 1.71        | 3.29        | 0.87        | 0.54        |
| LSTM                    | <b>2.15</b> | 4.97        | 0.93        | <b>0.75</b> | -           | -           | -           | -    | -           | -           | -           | -           |
| LSTM+social             | <b>2.15</b> | <b>4.95</b> | 0.93        | <b>0.75</b> | -           | -           | -           | -    | -           | -           | -           | -           |
| NN+map(prune)           | 3.38        | 7.62        | <b>0.99</b> | 0.86        | <b>2.11</b> | <b>4.36</b> | <b>0.97</b> | 0.68 | <b>1.68</b> | <b>3.19</b> | 0.94        | <b>0.52</b> |
| NN+map(prior) m-G,n-C   | 3.65        | 8.12        | 0.83        | 0.94        | 2.46        | 5.06        | <b>0.97</b> | 0.63 | 2.08        | 4.02        | <b>0.96</b> | 0.58        |
| NN+map(prior) l-G,n-C   | 3.65        | 8.12        | 0.83        | 0.94        | 3.01        | 6.43        | 0.95        | 0.80 | 2.6         | 5.32        | 0.92        | 0.75        |
| LSTM+map(prior) l-G,n-C | 2.92        | 6.45        | 0.98        | <b>0.75</b> | 2.31        | 4.85        | 0.97        | 0.71 | 2.08        | 4.19        | 0.95        | 0.67        |

Table 4: Motion Forecasting Errors for different number of predictions. minADE: Minimum Average Displacement Error, minFDE: Minimum Final Displacement Error, DAC: Drivable Area Compliance, MR: Miss Rate (with a threshold of 2 m). Please refer to Section 5.2 for definitions of these metrics (↓ indicates lower is better).

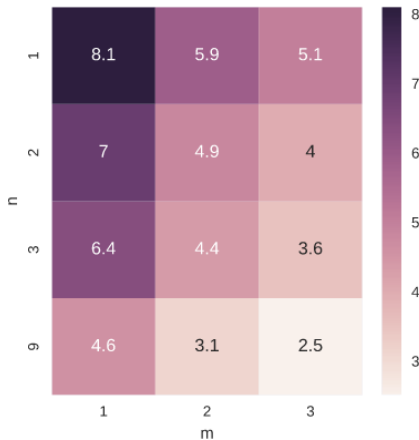


Figure 9: minFDE for  $NN+map(prior) m-G,n-C$  with different values of  $n$  (#Centerlines) and  $m$  (#Predictions along each centerline). There’s a trade-off between number of reference paths ( $n$ ) and number of predictions along each reference path ( $m$ ). Increasing  $n$  ensures that we are capturing different high level scenarios while increasing  $m$  makes sure we are capturing different velocity profiles along a given reference path. If the number of centerlines are enough, then for the same total number of predictions it is often better to make multiple predictions along fewer centerlines than to make 1 prediction along more centerlines.

## 6. Discussion

Argoverse represents two large-scale datasets for autonomous driving research. The Argoverse datasets are the first such datasets with 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 verify that map data can improve accu-

racy. We 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 our website [Argoverse.org](http://Argoverse.org).

**Acknowledgements.** We thank our Argo AI colleagues 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, Krathar 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 Govern, and Mike Webster. Odin: Team victortango’s entry in the darpa urban challenge. *J. Field Robot.*, 25(8):467–492, Aug. 2008.
- [4] Min Bai, Gellért Mátyus, Namdar Homayounfar, Shenlong Wang, Shrinidhi Kowshika Lakshmikanth, and Raquel Urtasun. Deep multi-sensor lane detection. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS 2018, Madrid, Spain, October 1-5, 2018*, pages 3102–3109, 2018.
- [5] Keni Bernardin and Rainer Stiefelhagen. Evaluating multiple object tracking performance: The clear mot metrics. *EURASIP J. Image and Video Processing*, 2008.
- [6] Holger Caesar, Varun Bankiti, Alex H. Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan,

- Giancarlo Baldan, and Oscar Beijbom. nuscenec: A multimodal dataset for autonomous driving. *arXiv preprint arXiv:1903.11027*, 2019.
- [7] 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.
- [8] 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.
- [9] 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.
- [10] 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.
- [11] Nachiket Deo and Mohan M Trivedi. Convolutional social pooling for vehicle trajectory prediction. *arXiv preprint arXiv:1805.06771*, 2018.
- [12] Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio Lopez, and Vladlen Koltun. Carla: An open urban driving simulator. *arXiv preprint arXiv:1711.03938*, 2017.
- [13] 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.
- [14] 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.
- [15] Alex Graves. Generating sequences with recurrent neural networks. *CoRR*, abs/1308.0850, 2013.
- [16] Junyao Guo, Unmesh Kurup, and Mohak Shah. Is it safe to drive? an overview of factors, challenges, and datasets for driveability assessment in autonomous driving. *IEEE Transactions on Intelligent Transportation Systems*, 2019.
- [17] 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.
- [18] 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.
- [19] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask R-CNN. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2017.
- [20] 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.
- [21] 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.
- [22] David Held, Jesse Levinson, and Sebastian Thrun. Precision tracking with sparse 3d and dense color 2d data. In *ICRA*, 2013.
- [23] 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.
- [24] Michael Himmelsbach and Hans-Joachim Wünsche. Lidar-based 3d object perception. In *Proceedings of 1st International Workshop on Cognition for Technical Systems*, 2008.
- [25] Namdar Homayounfar, Wei-Chiu Ma, Shrinidhi Kowshika Lakshmikanth, and Raquel Urtasun. Hierarchical recurrent attention networks for structured online maps. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [26] 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.
- [27] 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.
- [28] 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.
- [29] R. Kesten, M. Usman, J. Houston, T. Pandya, K. Nadhamuni, A. Ferreira, M. Yuan, B. Low, A. Jain, P. Ondruska, S. Omari, S. Shah, A. Kulkarni, A. Kazakova, C. Tao, L. Platinsky, W. Jiang, and V. Shet. Lyft level 5 av dataset 2019. <https://level5.lyft.com/dataset/>, 2019.
- [30] 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.
- [31] 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.
- [32] 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.
- [33] 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.
- [34] 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.
- [35] Justin Liang, Namdar Homayounfar, Wei-Chiu Ma, Shenglong Wang, and Raquel Urtasun. Convolutional recurrent network for road boundary extraction. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [36] Justin Liang and Raquel Urtasun. End-to-end deep structured models for drawing crosswalks. In *The European Conference on Computer Vision (ECCV)*, September 2018.
- [37] Ming Liang, Bin Yang, Yun Chen, Rui Hu, and Raquel Urtasun. Multi-task multi-sensor fusion for 3d object detection. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [38] Xingyu Liu, Charles R Qi, and Leonidas J Guibas. FlowNet3d: Learning scene flow in 3d point clouds. *arXiv preprint arXiv:1806.01411*, 2019.
- [39] 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.
- [40] 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.
- [41] 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.
- [42] 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.
- [43] 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.
- [44] Michael Montemerlo, Jan Becker, Suhrid Bhat, Hendrik Dahlkamp, Dmitri Dolgov, Scott Ettinger, Dirk Haehnel, Tim Hilden, Gabe Hoffmann, Burkhard Huhnke, 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.
- [45] 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.
- [46] 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.
- [47] 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.
- [48] 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.
- [49] 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.
- [50] Nicholas Rhinehart, Kris M Kitani, and Paul Vernaza. R2p2: A reparameterized pushforward policy for diverse, precise generative path forecasting. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 772–788, 2018.
- [51] Nicholas Rhinehart, Rowan McAllister, Kris Kitani, and Sergey Levine. Precog: Prediction conditioned on goals in visual multi-agent settings. *arXiv preprint arXiv:1905.01296*, 2019.
- [52] 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.
- [53] 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.
- [54] John Parr Snyder. *Map projections—A working manual*, volume 1395, page 61. US Government Printing Office, 1987.
- [55] 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.
- [56] 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.
- [57] 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 Ziglar. Tartan racing: A multi-modal approach to the darpa urban challenge. Technical report, Carnegie Mellon University, Pittsburgh, PA, April 2007.

- [58] Shenlong Wang, Min Bai, Gellert Mattyus, Hang Chu, Wenjie Luo, Bin Yang, Justin Liang, Joel Cheverie, 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.
- [59] Weiyue Wang, Ronald Yu, Qiangui Huang, and Ulrich Neumann. Sgpn: Similarity group proposal network for 3d point cloud instance segmentation. In *CVPR*, 2018.
- [60] 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.
- [61] Raymond A Yeh, Alexander G Schwing, Jonathan Huang, and Kevin Murphy. Diverse generation for multi-agent sports games. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4610–4619, 2019.
- [62] 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.

## Appendices

We present additional details about our map (Appendix A), our 3D tracking taxonomy of classes (Appendix B), our trajectory mining (Appendix C), and our 3D tracking algorithm (Appendix D).

### A. Supplemental Map Details

In this appendix, 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 m<sup>2</sup> in Miami and 286,104 m<sup>2</sup> in Pittsburgh.

#### A.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. In Pittsburgh, our city origin lies at 583710.0070 Easting, 4477259.9999 Northing in UTM Zone 17. In Miami, our city origin lies at 580560.0088 Easting, 2850959.9999 Northing in UTM Zone 17.

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 [54].

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 10 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 11 shows examples of centerlines, driveable area, and ground height projected onto a camera image.

#### A.2. Map API and Software Development Kit

The dataset’s rich maps are a novelty for autonomous driving datasets and we aim to make it easy to develop computer vision tools that leverage the map. Figure 13 outlines 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 egovehicle, 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 14).

These dense, pixel-level map renderings, similar to visualizations of instance-level or semantic segmentation [10], have recently been demonstrated to improve 3D perception and are relatively easy to use as an input to a convolutional network [60, 7].

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.

### B. 3D Tracking Taxonomy Details

Argoverse 3D Tracking version 1.1 contains 15 categories of objects. Examples of each object type are visualized in Table 5. Definitions of each of the 15 categories are given below.

- **Animal** A four-legged animal (primarily dogs or cats).
- **Bicycle** A non-motorized vehicle with 2 wheels that is propelled by human power pushing pedals in a circular motion.
- **Bicyclist** A person actively riding a bicycle.
- **Bus** A standard, city bus that makes frequent stops to embark or disembark passengers.



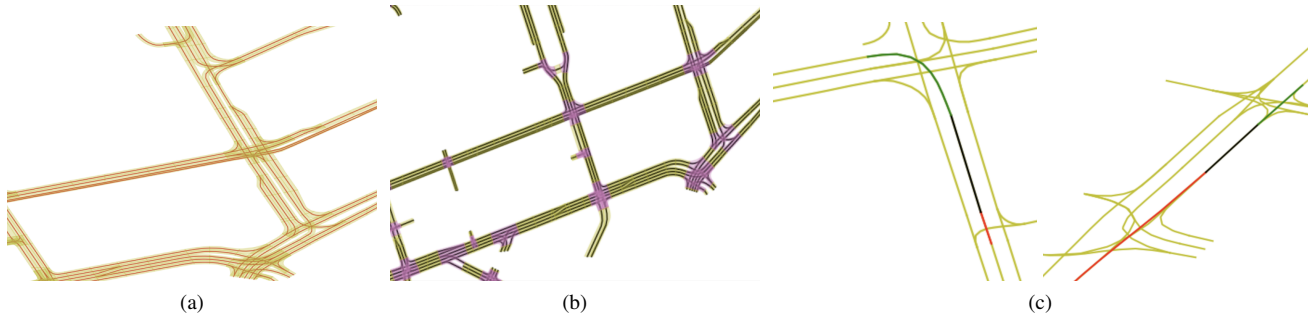


Figure 10: (a) Lane centerlines and hallucinated area are shown in red and yellow, respectively. We provide *lane* centerlines in our dataset because simple *road* centerline representations cannot handle the highly complicated nature of real world mapping, as shown above with divided roads. (b) We show 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.

- **Emergency vehicle** A vehicle with lights and sirens that, when active, gains right-of-way in all situations.
- **Large vehicle** Motorized vehicles with 4 or more wheels, larger than would fit in a standard garage.
- **Moped** A motorized vehicle with 2 wheels with an upright riding position with feet together.
- **Motorcycle** A motorized vehicle with 2 wheels where the rider straddles the engine.
- **Motorcyclist** A person actively riding a motorcycle or a moped.
- **On Road Obstacle** Static obstacles on driveable surface.
- **Other Mover** Movable objects on the road that don’t fall into other categories.
- **Pedestrian** A person that is not driving or riding in/on a vehicle.
- **Stroller** A push-cart with wheels meant to hold a baby or toddler.
- **Trailer** A non-motorized vehicle towed behind a motorized vehicle.
- **Vehicle** Motorized automobile- typically 4-wheels that could fit into a standard, personal garage.

## C. Supplemental Details on Motion Forecasting

In this appendix, we elaborate on some aspects of forecasting data mining and baselines.

### C.1. Motion Forecasting Data Mining Details

Here we describe our approach for mining data for trajectory forecasting. The scenarios that are 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. Furthermore, 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 overlap implies that the same track id can be available in 2 sequences, albeit with different trajectories.

### C.2. 2D Curvilinear Centerline Coordinate System

As discussed in Section 5, the baselines that use the map as a prior first transform the trajectory to a 2D Curvilinear Centerline coordinate system. In this section, we provide details about this new coordinate space. The centerline coordinate system has axes tangential and perpendicular to lane centerline. When a trajectory is transformed from the absolute map frame to the centerline frame, it makes the trajectory generalizable across different map locations and orientations. Figure 15 illustrates the transformation.

## D. Supplemental Tracking Details

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

### D.1. Tracker Implementation Details

Because of space constraints we provide details of our 3D tracker here instead of in the main manuscript. We do not claim any novelty for this “baseline” tracker, but it works reasonably

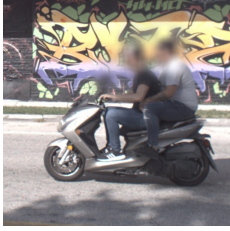
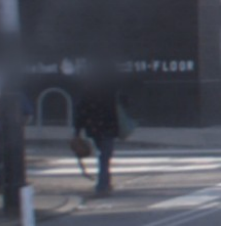
|  |  |  |  |   |
|--|--|--|--|---|
| ANIMAL   | BICYCLE<br>BICYCLIST   | BUS  | EMERGENCY<br>VEHICLE   | LARGE<br>VEHICLE  |
|   |   |   |  |  |
| MOPED  | MOTORCYCLE<br>MOTORCYCLIST   | ON ROAD<br>OBSTACLE  | OTHER MOVER  | PEDESTRIAN  |
|   |   |   |  |  |
| STROLLER   | TRAILER  | VEHICLE  |  |   |
|  |  |  |  |   |

Table 5: Argoverse 3D Tracking version 1.1 categories.

well, especially with map information available (e.g. driveable area, ground height, and lane information). More recent attempts in 3D tracking and detection include the baseline introduced in H3D [47], with VoxelNet [62] detection and an Unscented Kalman Filter [28, 27] for tracking. The NuScenes dataset [6] uses PointPillars [31] for a 3D bounding box detection baseline. Apollo-Car3D [55] implements 2D to 3D car pose estimation baselines based on 3D-RCNN[30] and DeepMANTA [8]. However, we do not compare our baseline tracker against these methods.

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, we use geometrically cluster the remaining 3D LiDAR point cloud into separate objects according to density, using DBSCAN [13]. Then we use Mask R-CNN [19] to obtain object masks in pixel space and discard any LiDAR clusters whose image projection does not fall within a mask.

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 [21]. We can avoid many such segmentation failures

with the high precision of our Mask R-CNN network<sup>2</sup>. We also alleviate the need for an object’s full, pixel-colored 3D shape during tracking, as others have suggested [23, 22]. We prefer density-based clustering to connected components clustering in a 2D occupancy grid [24, 34] 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 driveable region defined by the map. We also perform ground removal using either the ground height map or plane-fitting.

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.

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



(a) Lane geometry and connectivity



(b) Driveable area



(c) Ground height

Figure 11: Examples of centerlines, driveable area, and ground height projected onto a camera image.

**3. Tracking.** We use ICP (Iterative Closest Point) from the Point Cloud Library [53] to estimate a relative transformation between corresponding point segments for each track. Then we apply a Kalman Filter (KF) [22] 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.

## D.2. Tracking Evaluation Metrics

We use standard evaluation metrics commonly used for multiple object trackers (MOT) [43, 5]. The MOT metric relies on centroid distance as distance measure.

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

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

where  $FN_t$ ,  $FP_t$ ,  $ID_{sw}$ ,  $GT$  denote the number of false negatives, false positives, number of ID switches, and ground truth objects. We report MOTA as percentages.

- **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):** Ratio of trajectories tracked more than 80% of its lifetime.
- **ML (Mostly Lost):** Ratio of trajectories tracked for less than 20% of object lifetime, over the entire object 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".

## D.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 operate on *pixels* which could be any real world size, whereas in 3D we have absolute lengths. When using IoU as a TP/FP threshold, small objects tend to be penalized unfairly. For example, for pairs of bounding boxes with the same distance between centroids (e.g. a pedestrian that is tracked with 0.5 m error vs a car that is tracked with 0.5 m error), the larger objects will have higher IoU (see Figure 16). To track pedestrians with the same IoU as buses requires orders of magnitude more positional precision.

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. Still, the centroids of such objects can be estimated. For these reasons, we use the absolute distance between centroids as the basis for classifying correct vs. incorrect associations between tracked and ground truth objects.





Figure 12: **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> | Uses rasterized driveable area ROI to decimate LiDAR point cloud to only ROI points.  |
| <code>remove_ground_surface</code>            | Removes all 3D points within 30 cm of the ground surface.   |
| <code>get_ground_height_at_xy</code>          | Gets ground height at provided (x,y) coordinates.   |
| <code>render_local_map_bev_cv2</code>         | Renders a Bird's Eye View (BEV) in OpenCV.  |
| <code>render_local_map_bev_mpl</code>         | Renders a Bird's Eye View (BEV) in Matplotlib.  |
| <code>get_nearest_centerline</code>           | Retrieves nearest lane centerline polyline.   |
| <code>get_lane_direction</code>               | Retrieves most probable tangent vector $\in \mathbb{R}^2$ to lane centerline.   |
| <code>get_semantic_label_of_lane</code>       | Provides 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>          | Gets all lane IDs within a Manhattan distance search radius in the <i>xy</i> plane.   |
| <code>get_lane_segment_predecessor_ids</code> | Retrieves all lane IDs with an incoming edge into the query lane segment in the semantic graph.   |
| <code>get_lane_segment_successor_ids</code>   | Retrieves all lane IDs with an outgoing edge from the query lane segment.   |
| <code>get_lane_segment_adjacent_ids</code>    | Retrieves all lane segment IDs of that serve as left/right neighbors to the query lane segment.   |
| <code>get_lane_segment_centerline</code>      | Retrieves polyline coordinates of query lane segment ID.  |
| <code>get_lane_segment_polygon</code>         | Hallucinates a lane polygon based around a centerline using avg. lane width.  |
| <code>get_lane_segments_containing_xy</code>  | Uses a "point-in-polygon" test to find lane IDs whose hallucinated lane polygons contain this (x,y) query point.                                |

Figure 13: Example Python functions in the Argoverse map API.



Figure 14: A scene with non-planar ground surface. The colored LiDAR returns have been classified as belonging to the 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.

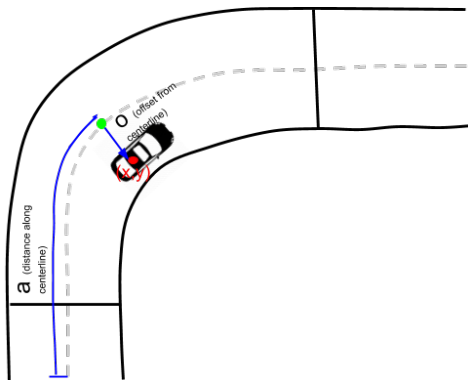
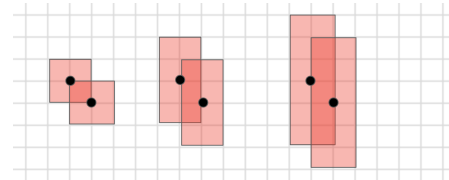
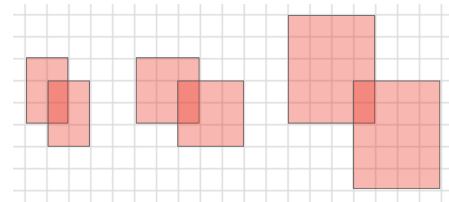


Figure 15: Converting a coordinate  $(x, y)$  from the absolute map frame to the 2D curvilinear centerline coordinate frame  $(a, o)$  where a vehicle's position is described by a distance  $a$  along a centerline and an offset  $o$  from the centerline.

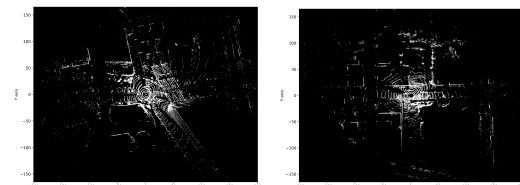


(a) Fixed inter-centroid distance of  $\sqrt{2}$  m.



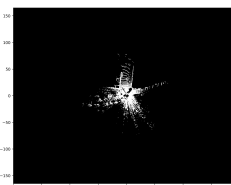
(b) Fixed 2D intersection area of  $2 m^2$ .

Figure 16: 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.



(a) Argoverse LiDAR

(b) Argoverse LiDAR



(c) KITTI LiDAR



(d) nuScenes LiDAR

Figure 17: 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 up to twice the range of the sensors used to collect the KITTI or nuScenes datasets, allowing us to observe more objects in each scene.