

# DeepAutoTrack (DAT): Vehicle Trajectory Prediction for Autonomous Driving

Anonymous CVPR submission

Paper ID 3790

## Abstract

Vision-based deep neural networks are crucial components for autonomous driving in terms of visual understanding and decision making. These models need to be temporally consistent and visually interpretable. Most recent works focus on using static images for understanding visual semantics, ignoring the temporal consistency of driving conditions; or adopt end-to-end architectures for learning the pilot networks, providing limited interpretability. In this paper, we raise the question of “can we predict other traffic participant’s future trajectory given previous egomotion visual input?”. The proposed system stems from the issue of human driver reaction time in the autonomous driving context. Our dual-staged model firstly applies modern convolutional object detectors for spotting traffic participants (i.e., cars in this paper) and robustly tracks the targets of interest. Then, a novel SEG-LSTM network is incorporated to fuse the multiple-streams from past frames and then predict targets’ future trajectories. We demonstrate the feasibility of predicting future trajectories of vehicles for assisting mediated perception. We also show the effectiveness of combining various levels of abstractions(e.g., scene semantics, detection bounding boxes) further boost the prediction accuracy across diverse traffic conditions.

## 1. Introduction

Currently, there are two major paradigms for autonomous driving systems built upon vision-based input [6]: mediated perception approaches that firstly explain the vision input and then parse the scene to make a driving policy (usually by a controller with if-then-else rules); and behaviour reflex approaches that directly map the vision input to a driving policy by a regressor. In this paper, in the framework of first paradigm, we try to tackle the problem of predicting future trajectories of vehicle given its past information.

Being able to predict other traffic participants’ future trajectories is important because it can help to prevent self-driving car from running into other cars. We need to know

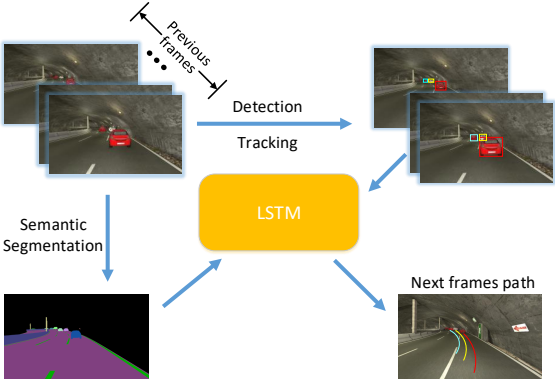


Figure 1: We propose a system to predict future trajectories of vehicles. Given historical information, obtained by detection, tracking and semantic segmentation, a multi-stream recurrent neural network is trained to generate future trajectory.

not just where other cars are, as in the localization case, but also how fast they are moving in order to avoid collisions. For a human driver, reaction time is a crucial factor that includes recognizing the light has changed, deciding to continue or brake, and whether to stop engaging the brake (remove foot from accelerator and apply brake). Accident reconstruction specialists commonly use 1.5 seconds [35]. Therefore, the ability to predict traffic participants’ future trajectories will greatly benefit the mediated perception approach’s driving policy decision making.

Nonetheless, vehicles trajectory prediction is a challenging problem: by dissecting visual scene understanding individually (e.g., object detection, semantic segmentation, instance segmentation), there are still many open challenges for computer vision community. In addition, the scene parsing results contain quite some spurious information. Existing methods for scene parsing also mainly focus on static images. However, autonomous driving is intrinsically a dynamic problem. Human drivers make decisions by a sequence of input frames instead of one static frame. Therefore, temporal consistency among continuous frames should also be taken into consideration when designing the interpretable system.

Vehicle trajectory prediction is also very related to visual tracking. In traditional vision based tracking problems (e.g., [48, 37]) where target objects go through rather sporadic movements (usually the object trajectory is generated by artificial movement in order to test the robustness of the tracker), the prediction of the target’s future trajectory is simply a moot problem. However, traffic participants generally exhibit regular trajectories by obeying a large number of common sense rulesp (e.g., cars driving on the road, pedestrians walking on the sidewalks). Given law-abiding traffic participants, human drivers subconsciously project visual target’s future trajectory. Nonetheless, due to the missing bridge between the traditional visual tracking community and the current autonomous driving research community, there is a lack of proper metrics for evaluating the temporal prediction result.

We address this issue through building a system that takes advantage of multiple streams with various levels of abstraction. In our approach, there are three key components: traffic participants detection, instance tracking and future trajectory prediction as shown in Fig. 1. Our main contributions are as follows:

- We construct a time-series dataset for vehicle trajectory prediction based on the SYNTHIA dataset [41]<sup>1</sup> and verify the feasibility of predicting vehicles’ future trajectories for assisting mediated perception.
- We take advantage of various levels of abstractions from multiple streams via recurrent network based temporal models and demonstrate the effectiveness of using “auxiliary information”(e.g., scene semantics) further boost the prediction accuracy.
- Preliminary results by incorporating “social” factor demonstrate the potential of learning common sense rules with neighboring traffic participants in complex real world environment.

## 2. Related Work

**Traffic participants detection** Traffic participants detection is one key component of an autonomous driving system. Vehicle detection has been well studied for static highway mounted cameras: [22] present a concise overview of image processing methods and analysis tools before the “deep learning eras”. A lot of progress has been made in recent years on object detection due to the use of convolutional neural networks (CNNs). Modern object detectors [20, 19, 40, 24, 23, 39, 33] are now good enough to be deployed in a consumer products. Faster R-CNN [40] proposes a region proposal network (RPN) and integrates RPN and Fast R-CNN [19] to build an end-to-end architecture. SSD [33] similarly uses a series of boxes and directly

<sup>1</sup>Dataset will be released upon paper publication.

gives the probability for each object class.

Typical algorithms output bounding boxes on detected objects for a static image. And the deep neural nets can be transferred and fine-tuned to another set of classes for traffic participants we actually care about in the context of autonomous driving.

**Instance association** Instance association between frames is a prerequisite to solve the proposed intrinsically temporal problems. 3D scene flow is estimated in [3] by exploiting recognition to overcome the challenging scenarios in the presence of large displacement or local ambiguities. However, a total runtime of 10 minutes per frame is computationally prohibitive for the real-time deployment. To accomplish the task of instance association between frames, we resort to the well studied field of visual tracking. Discriminative Correlation Filters (DCF) [26] have demonstrated excellent performance for high-speed generic visual object tracking. Built upon their seminal work, there has been a plethora of recent improvements [29, 34, 12, 25, 46, 10] relying on convolutional neural network (CNN) pre-trained on ImageNet as a feature extractor for visual tracking.

Most of the modern trackers focus on the problem of “class-agnostic” generic object tracking. In order to adapt to temporal changes, a continuous learning strategy is applied, where the model is updated rigorously in every frame. Such update is excessive and sensitive to sudden changed caused by, e.g., scale variations, deformations, and out-of-plane rotations. This excessive update strategy cause both lower frame-rates and degradation of robustness due to over-fitting to the recent frames [10].

**Activity forecasting** Activity forecasting is the task of predicting the motion or the action to be carried out by objects in a video. A body of work [5, 8, 15] explores various hierarchical models in face of different quality of data representation. Recently recurrent neural networks (RNN) and their variants including Long Short Term Memory (LSTM) [28] and Gated Recurrent Units (GRU) [7] have proven to be very successful for sequence prediction tasks. An action-conditioned video prediction model [17] that explicitly models pixel motion, by predicting a distribution over pixel motion from previous frames is proposed to learn from unsupervised video data. A curriculum learning strategy is proposed in [4] to gently change the training process from a fully guided scheme using the true previous token, towards a less guided scheme which mostly uses the generated token instead. A combination of a fully-convolutional network and an LSTM is proposed in [49] to learn from large-scale vehicle action data and is trained as a pilot network in the context of behaviour reflex approaches. Their goal is to predict future egomotion including multimodal discrete and continuous driving behaviors. Inverse turning radius is predicted by a LSTM network in [32]: a

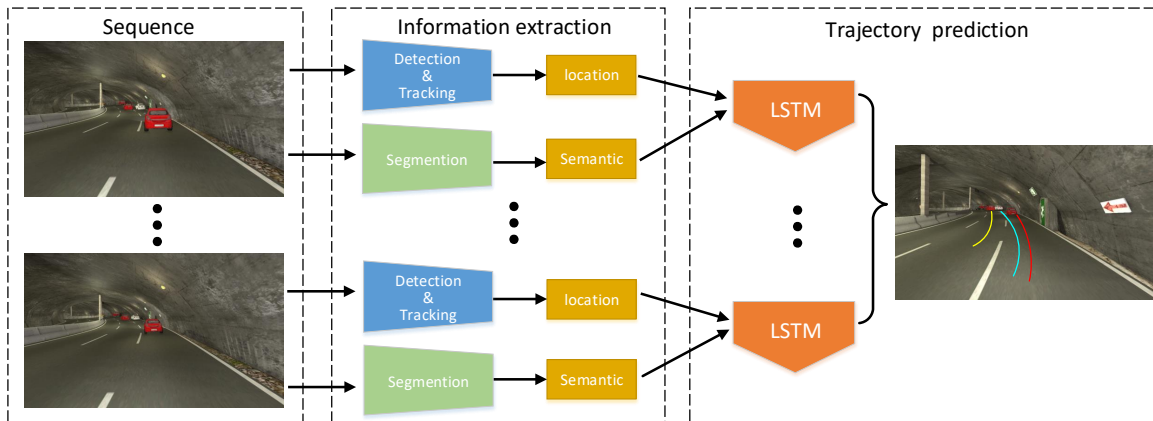


Figure 2: The framework for the proposed vehicle future trajectory prediction system. The input is the sequential data (RGB and/or depth maps) up to the current time step obtained from the host car. Target vehicles’ sequential bounding boxes are extracted by detection & tracking module. Along with scene semantics as an auxiliary information, multi-stream inputs are fed into an LSTM to predict future trajectories taking advantages of various levels of abstractions.

heat map of attention is generated at each time step conditioned on the previous hidden states as a more succinct visual explanations.

More recent works also attempt to take the “social factor” into consideration. Social-LSTM [2] focus on modeling dynamics crowd interactions for path prediction by introducing a “social tensor” from a spatial-temporal neighborhood. Based on the Social-LSTM, an attention mechanism is introduced in [45] built upon the high-level spatio-temporal graphs proposed in [31].

**Datasets for autonomous driving** Vision-based semantic segmentation in urban scene is a key functionality for autonomous driving. KITTI benchmark suite [18] provides a large amount of images of urban scene from Karlsruhe, Germany, with ground truth data for odometry, object detection, tracking, scene flow and stereo benchmarks. However, a limited 430 labelled images are provided for semantic segmentation. More recently, larger projects are constructed: Cityscapes dataset [9] which consists of a collection of images acquired in 50 cities around Germany, Switzerland and France in different seasons, and having 5,000 images with fine annotations and 20,000 with coarse annotations over a total of 30 classes. Comma.ai [42] provides 7 and a quarter hours of largely highway driving. Udacity [1] dataset includes 65,000 labels across 9,423 frames at full resolution of  $1920 \times 1200$  at 2Hz. The use of synthetic data has increased considerably in recent years within computer vision community. CARLA [14] and Sim4CV [36] are two recent endeavours that take advantage of simulated environment to support development, training and validation of autonomous urban driving systems. A synthetic dataset named SYNTHIA [37, 27] is collected for investigating how useful synthetic images can be for semantic segmentation. Syn-

thetic dataset has an obvious advantage: a fine annotated image in Cityscapes dataset requires on average 1.5 hours which is very labor intensive. Thus, the cost of scaling large project would require a prohibitive economic investment in order to capture images from a larger variety of countries, in different seasons and different traffic conditions.

### 3. Proposed Algorithm

We first formalized the proposed problem and describe the overall framework as shown in Fig. 2.

#### 3.1. Problem formulation

We propose to learn a generic approach from history information and formulate the problem as predicting future traffic participants’ trajectories. Formally, our problem can be defined as a mapping  $\mathcal{M}$  from historical trajectory  $\mathcal{J}_H = \{j_p^t\}, t = 1, \dots, T$  with auxiliary information  $\mathcal{A} = \{a_p^t\}, t = 1, \dots, T$  for traffic participant  $p$  to the future trajectory  $\mathcal{J}_F = \{j_p^t\}, t = T + 1, \dots, T_F$ :

$$\mathcal{J}_F \leftarrow \mathcal{M}_p(j, a) : \mathcal{J}_H \times \mathcal{A} \quad (1)$$

where the traffic participants  $p$  can be vehicles, cyclists, pedestrians, etc. In this paper, we focus on vehicle as the sole traffic participant category and ignore the subscript  $p$  from now on.  $T_F$  is the length of future trajectory to be predicted (corresponding to the total driver reaction time). Historical and future trajectories are defined in a 3D occupancy grid:

$$\mathcal{J}_H, \mathcal{J}_F \in \mathbb{R}^6 = \{x, y, w, h, d_{min}, d_{max}\} \quad (2)$$

where  $\{x, y, w, h\}$  define target’s 2D boundingbox in image pixel space and  $\{d_{min}, d_{max}\}$  define the minimum

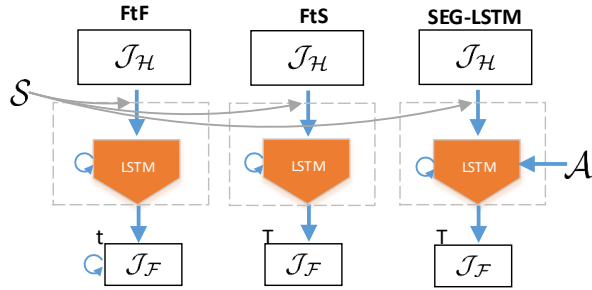


Figure 3: Illustrations of three temporal networks.  $\mathcal{J}_H$  represents historical trajectories,  $\mathcal{J}_F$  is the  $T$  time step future trajectories to be predicted. Left: Frame-to-Frame (FtF); Middle: Frame-to-Sequence (FtS); Right: SEG-LSTM, with auxiliary information as input  $A$ . The social tensor  $S$  can be inserted into existing recurrent networks.

and maximum relative distance between host car and the tracking target (acquired from the depth camera). Given RGB cameras' focal length  $f$ , the pixel space 2D boundingbox can be projected into 3D real-world coordinates  $\{x_r, y_r, w_r, h_r\}$  as:

$$x_r = \frac{x}{f}d, y_r = \frac{y}{f}d \quad (3)$$

$$w_r = \frac{w}{f}d, h_r = \frac{h}{f}d \quad (4)$$

where  $d$  is distance between the center of the camera lens and the center of target. We use the minimum relative distance  $d_{min}$  as a close approximation.

### 3.2. Sequence modeling

For predicting the vehicle's future trajectory, we compare three temporal models based on LSTM cells. A "social pooling tensor" that implicitly reasons the motion of neighboring traffic participants can be inserted into existing recurrent network architectures without any extra training supervision or modification to the optimization process as shown in Fig. 3. The first two models take sequential 3D occupancy as input. In the third model, we explore using semantic segmentation as a side "auxiliary" information  $A$  to better guide the trajectory prediction.

**Frame-to-Frame (FtF) model:** is akin to sequence generation model from Graves [21]. At test-time  $t$ , it samples current predicted 3D occupancy as feedback to the model for predicting next frame occupancy grid. The model is frame-to-frame Markovian and the next frame trajectory  $j_{t+1}$  can be estimated as:

$$j_{t+1} \leftarrow \mathcal{M}(j_t, j_H) : \mathcal{J}_H \quad (5)$$

**Frame-to-Sequence (FtS) model:** differs from the FtF

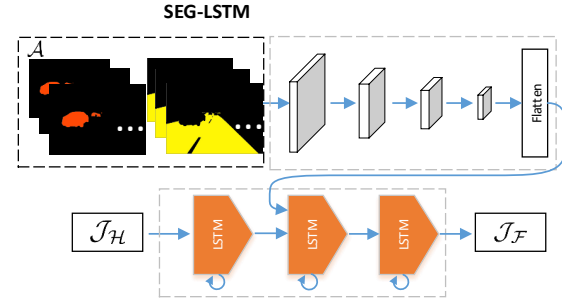


Figure 4: Diagram of SEG-LSTM model: semantic inputs are one-hot encoded images with channel number  $C$  as the number of semantic categories (we used only "car" and "road" as the two semantics under the context of high-way driving). The higher level semantics are extracted through a four-layer CNN. A joint representation is produced by the auxiliary semantic information and the coordinates trajectory information.

model in that instead of sampling frame by frame for predicting future trajectories, FtS model directly predicts a trajectory sequence of length. As it pointed out in [4] that FtF training paradigm, predicting one token at a time, conditioned on the state and the previous correct token, is different from how we actually use them during inference and thus is prone to the accumulation of errors along the decision paths. Therefore, a curriculum learning approach is proposed to slowly change the training objective from an easy task, where the previous token is known, to a realistic one, where it is provided by the model itself. The aforementioned learning strategy applies to inference for time-varying sequence where an  $\langle \text{EOS} \rangle$  token is needed to signify the end of the sequence. For our continuous prediction task, we stipulate the length of the prediction is predefined as  $T_F$ . Hence we have the assumption that current state is complete, in the sense it contains all historical information required for predicting the future trajectory of the vehicle. The whole trajectory  $j_{t+1, \dots, T_F}$  to be predicted can be estimated as:

$$j_{t+1, \dots, T_F} \leftarrow \mathcal{M}(j_H) : \mathcal{J}_H \quad (6)$$

The whole sequence inference paradigm is more stable than the first FtF approach because of the elimination of the dependency of previous token. Such dependency will introduce the discrepancy between how the model is used at training and inference. Inference with the whole sequence will generate smoother predictions.

**SEG-LSTM model:** takes advantages of scene semantics from historical frames as auxiliary information  $A$ . The second LSTM layer fuses the output from the bottom-most LSTM that encodes coordinates embedding with higher level semantic information. With auxiliary  $A$ , the trajectory



$j_{t+1, \dots, T_F}$  to be predicted can be estimated as:

$$j_{t+1, \dots, T_F} \leftarrow \mathcal{M}(j_H, a) : \mathcal{J}_H \times \mathcal{A} \quad (7)$$

Semantic inputs are one-hot encoded images with channel number  $\mathcal{C}$  as the number of semantic categories as in Fig. 4. The higher level semantic is extracted through a four-layer CNN. A joint representation is produced by the semantic auxiliary information and the coordinates trajectory information. As described above, the 3D occupancy input has only six dimensions:  $[x, y, w, h, d_{min}, d_{max}]$  encoding the coordinates information. So the bottom-most LSTM not only can learn the temporal information between time-varying trajectories, but also can balance input dimensions of both coordinates information and the semantic information.

Note that this model differs from the way semantic information is used in [49]: instead of using semantic segmentation as a “privilege information” which is missing at test time, SEG-LSTM directly takes historical scene semantics as input and embeds them to a hidden-state tensor by a convolutional neural network.

**Social tensor:** was firstly introduced in [2] by observing the “common sense” rules and social conventions from pedestrians in crowded public space. Naturally, in the context of autonomous driving, such observation is also complied: when there are multiple cars occupying on-going lanes, the social interactions between traffic participants confine the space for their possible trajectories. The social interactions can be encoded implicitly by using a social pooling layer: at every time-step, the LSTM cell receives pooled hidden-state information from the LSTM cells of neighbors. While pooling the information, the spatial information is preserved by grid based pooling.

The hidden state  $h_i^t$  of the LSTM at time  $t$  captures the latent representation of the  $i^{th}$  target in the scene. This representation can be shared with neighbors by building a social hidden tensor  $s_t^i$ :

$$s_t^i = \sum_{j \in \mathcal{N}_i} \mathbf{1}_{m,n}[x_t^j - x_t^i, y_t^j - y_t^i] \cdot h_{t-1}^j \quad (8)$$

where  $\mathbf{1}_{m,n}[x_t^j - x_t^i, y_t^j - y_t^i]$  is an indicator function to check if  $(x, y)$  is in the  $(m, n)$  cell of the grid, and  $\mathcal{N}_i$  is the set of neighbors corresponding to target  $i$ . The trajectory to be predicted can be estimated as:

$$j_{t+1, \dots, T_F} \leftarrow \mathcal{M}(j_H, s_H^i) : \mathcal{J}_H \times \mathcal{S} \quad (9)$$

This social pooling layer does not introduce any additional parameters and can be inserted into the hidden layer of LSTM. An important distinction from the vanilla LSTM is that the hidden states of multiple LSTMs are coupled by the social pooling layer. At every time-step, the updates are

jointly back-propagated through multiple LSTMs in a scene as a neighboring batch update scheme.

## 4. Experiments

We first provide the data collection procedure for vehicle trajectory prediction. And then we present the implementation details with performance analysis of individual modules. Quantitatively examination of SEG-LSTM architecture is also presented showing the advantages of using semantics directly as input can better assist the neural net’s trajectory prediction.

### 4.1. Dataset collection

Currently, most public datasets with semantic labeling are single frame, static images. Continuous video streams are required for the purpose of car trajectory prediction. We collect a car trajectory prediction dataset based on the SYNTHIA [41] dataset. SYNTHIA dataset is a large corpus of synthetic images originally collected for the purpose of semantic segmentation of urban scenes generated by rendering a virtual city created with the Unity development platform. The potential of this virtual world includes extension capabilities: new parts of the cities and road conditions can be easily generated by adding different setups. The major features of the collected dataset include: scene diversity (European style town, modern city, highway and green areas), variety of dynamic objects (cars, pedestrians and cyclists), multiple seasons (dedicated themes for winter, fall, spring and summer), lighting conditions and weather (dynamic lights and shadows, several day-time modes, rain mode and night mode). There are more than 200,000 HD ( $760 \times 1280$ ) photo-realistic frames from video streams. Frames are acquired from multiple view-points (up to eight views per location), and each of the frames also contains an associated depth map. In our experiment, we used only left front camera view data.

The focus of this paper is on car trajectory prediction on highways which corresponds to sequence number 1 in the SYNTHIA dataset. We are especially concerned about the cars that are relatively close to the driver. In order to decide quantitatively which vehicles are close to the host car, we define pixel occupancy ratio (*POR*) as the ratio between the total number of pixels of the tracking vehicle and the total number of pixels of the camera input. As in accordance with [6], we set the reliable car perception as 30 meters so as to guarantee satisfactory control quality when the speed of the host car does not exceed 72km/h. Since depth information is also provided alongside in the SYNTHIA dataset, we back-project the segmented instance region onto the depth maps and estimate that the *POR* of 0.2% and 0.1% correspond to roughly 30 meters and 100 meters of relative distance between the host car and the tracked vehicle. Hence, we start tracking the target when the *POR* of

**Ground truth collection for car trajectory prediction**

**Step 1** → acquiring traffic participants from ground truth annotations with only “car” as instances.

**Step 2** → acquiring frame-based instance information:

if  $POR > 0.2\%$  then

⇒ start tracking instance;

⇒ collecting 3D tracking boundingbox:

$[x, y, w, h, d_{min}, d_{max}]$ ;

⇒ collecting semantic labeling image  $I_{seg}$ ;

⇒ collecting 2D car pose from RGB input  $I_{car}$ ;

else if  $POR < 0.1\%$  then

⇒ stop tracking.

end if

**Step 3** → splitting sequences into tracklets of 23 frames with stepsize 1:

⇒ first 15 frames (5Hz, 3 sec) as training input;

⇒ next 8 frames (1.6 sec) as the held-up future trajectory to be predicted.

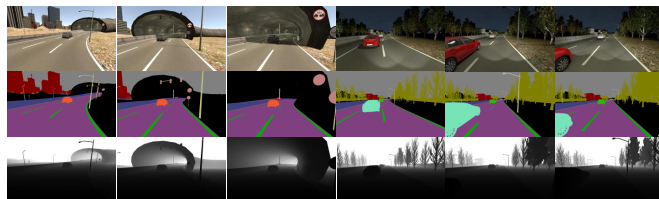


Figure 5: Examples of collected dataset for vehicle trajectory prediction. From top to bottom: RGB, semantic labeling, depth maps for three continuous frames under two dynamic light conditions.

the detected vehicle is larger than 0.2% so as to suffice the safety distance for driver to make timely reaction and stop tracking the vehicle when the POR is smaller than 0.1% indicating the target car is too far to influence driving policy. We follow the prediction time allocation as in [49] that the historical time-span is 3 seconds (corresponds to 15 frames with 5 Hz frame rate in the SYNTHIA dataset). However, in contrast with [49] where only the next 1/3rd of a second is predicted, we strive to predict the next 1.6 second’s future trajectory as mentioned previously that it’s the reaction time accident reconstruction specialists commonly used [35].

We follow the traditional paradigm for visual tracking to collect target’s bounding box:  $[x, y, w, h]$  of each frame. Moreover, the minimum relative distance  $d_{min}$  and maximum relative distance  $d_{max}$  can also be acquired from the depth camera. The train/valid/test are split as 80%/10%/10% of the total tracklets and the total numbers are shown in Tab. 1. Some sample frames with semantic labels and depth information is shown in Fig. 5.

	#train	#valid	#test
tracklets	10400	1300	1280
car detection	7832	976	986
semantic segmentation	7841	977	987

Table 1: Dataset statistics

## 4.2. Implementation Details

### 4.2.1 Traffic participants detection

To reliably detect traffic participants, we compared two state-of-the-art approaches: SSD [33] and Faster-RCNN [40]. We use the pretrained network on the PASCAL VOC detection dataset [16] which has 20 classes and fine-tune the network on the SYNTHIA with two classes: cars *v.s.* background. A comprehensive survey of trade-offs for modern convolutional object detectors is presented in [30] and we refer keen readers to the aforementioned paper for a more complete comparison to achieve the right speed/memory/accuracy balance for a given application and platform.

As is also pointed out in the paper [30] that SSD, albeit is less advantageous in detecting small objects, it’s very competitive for detecting larger objects. In this paper, we are mostly concerned about cars that are close to the driver. Hence, given 0.1% POR as the cut-off threshold, the presented problem favors detectors that are robust in detecting larger objects. Tab. 5 verifies that SSD is indeed more competitive when objects of interests are large in our problem set. We also compare the influence of the cut-off threshold for various confident scores and non-maximum suppression (NMS) thresholds. It can be seen that both meta-architectures are robust to cut-off confident scores. Allowing larger NMS threshold enables the detector to have slightly higher recall, it comes at a cost of multiple redundant overlapping detections and much lower f-score overall. In the following experiment, we adopt SSD with a confident score of 0.5 as the cut-off threshold and 0.45 as the default jaccard overlap for NMS.

### 4.2.2 Instance association via tracking

To associate traffic participants across different frames (*i.e.*, instance association), we creatively combine detecting with tracking. For this task, we adopt the real-time scale adaptive tracker fDSST [11] that achieves 50 FPS with scale estimation. For every newly detected target whose POR is larger than the predefined threshold of 0.2%, we initiate the tracker with a unique target ID.

The tracker’s target model is updated every time when there is a detection whose jaccard overlap is larger than 0.3. We argue that force update of tracker model via the detection template is essential when the target undergo rapid

out of plane rotations. Trackers [38, 13, 46, 10] focus on the “class-agnostic” settings. In the context of autonomous driving, however, we know in advance the classes of object of interests for tracking. Therefore, with this extra source of information, we propose to use the neural nets object detection result as a more robust template to update the trackers model to overcome the problem of model drifting. Hence the tracker is served solely as instance associator for the traffic participant. In return, the instance association scheme helps to alleviate the problem of false positive detection (*c.f.* Fig. 8) from the detection result. The target’s tracker will terminate under either of the following two conditions: if the target’s POR is smaller than 0.1% (*i.e.*, the vehicle is too far way from the host car); or if there are more than 5 frames of detection absence (*i.e.*, the target could be one false positive detection from the vehicle detector).

### 4.2.3 Trajectory prediction

For the task of sequential prediction, both FtF and FtS share the same base network structure: a 3-layer recurrent net with 100, 300, 300 LSTM cells. For SEG-LSTM, since the focus of the scene context is high-way driving, auxiliary information is the scene semantics with number of channels  $\mathcal{C}$  sets as 2 (“car” and “road”). We adopt the state-of-the-art semantic segmentation networks: Dilated Residual Networks (DRN) [50] and retrain on the SYNTHIA corpus. DRN achieves 68% and 82% pixel mAP on the car and road categories. The input semantics  $\mathcal{A}$  are one-hot encoded in the dimension of  $95 \times 160 \times 2$ . The CNN for learning high level semantics is composed of 4 convolutional blocks. Each block includes a convolutional layer, a ReLU activation layer, a maxpooling layer and a dropout layer (as we have observed overfitting during network training). Each convolutional layer has four channels with kernel of size 3.

For the social tensor configuration, we incorporate all vehicles in the same time-step as their neighbors. In contrast with [2] where a spatial constraint is required in order to be counted as a neighbor, all the vehicles in the same scene are relevant for the social interactions under our high-way driving context. Grid size is set as 4 and the embedding size is 64 so the social tensor  $\mathcal{S}$  is a  $4 \times 4 \times 64$  matrix. The learning rate is set to 0.001 with RMS-prop [44] with a gradient clipping of 3 to avoid divergence.

### 4.3. Evaluation

**Baseline: Kalman filter**, also known as linear quadratic estimation, is a popular technique for estimating the state of a system [43]. For the task of trajectory prediction, Kalman filters estimates a continuous state and gives a uni-modal distribution. The Kalman filter represents all distributions of the Gaussians and iterates two states: (1) measurement updates; (2) motion updates. In the baseline, we set the ini-

Configuration	mean coverage	center error
2D	(%)	(pixel)
Kalman Filter	54.31	24.84
<b>FtF</b>	66.47	23.43
<b>FtS</b>	73.93	15.32
<b>SEG-LSTM</b>	<b>77.53</b>	<b>12.63</b>

Table 2: Results of 2D performance evaluation (pixel as unit) on mean coverage (higher is better) and center error (lower is better).

Configuration	mean coverage	center error
3D	(%)	(meter)
<b>FtF</b>	35.67	0.7122
<b>FtS</b>	46.69	0.3686
<b>SEG-LSTM</b>	<b>48.49</b>	<b>0.3298</b>

Table 3: Results of 3D performance evaluation (centimeter as unit) on mean coverage (higher is better) and center error (lower is better).

tial uncertainly covariance as  $1e4$ , measurement noise and motion noise are all set with unit value.

**2D Evaluation Methodology:** Following the evaluation strategy of [47], all trajectories are compared using two measures: precision and success. Precision is measured as the distance between the centers of the ground truth bounding box and the corresponding tracker generated bounding box. The precision plot shows the percentage of tracker bounding boxed within a given threshold distance in pixels of the ground truth. To rank the prediction performance, the conventional threshold of 20 pixels ( $P20$ ) is adopted. Success is measured as the intersection over union of pixels. The success plot shows the percentage of tracker bounding boxes whose overlap score is larger than a given threshold and the trackers are ranked according to the Area Under Curve ( $AUC$ ) criteria.

**3D Evaluation Methodology:** Given depth map and the parameter of camera focal lens, we can formalize the problem into a 3D occupancy grid estimation. The real world coordinates  $[x_r, y_r, w_r, h_r]$  can be obtained according to Eq. 3 4. The 3D jaccard index can be similarly formalized as the volume intersection over union:  $\frac{V_g \cap V_p}{V_g \cup V_p}$ , where  $V_g, V_p$  are the ground truth and predicted space occupancy. 3D center error correspondingly is defined as the mean square root in 3D space.

**Learning with auxiliary information:** The 2D and 3D evaluation results in Fig. 6 and Tab. 2 show that our proposed system consistently outperforms the baseline Kalman filter. Kalman filter is intrinsically a linear model and it’s difficult to handle highly nonlinear driving trajectories in a relative coordinates. The auxiliary learning with semantic information in SEG-LSTM also consistently outperforms



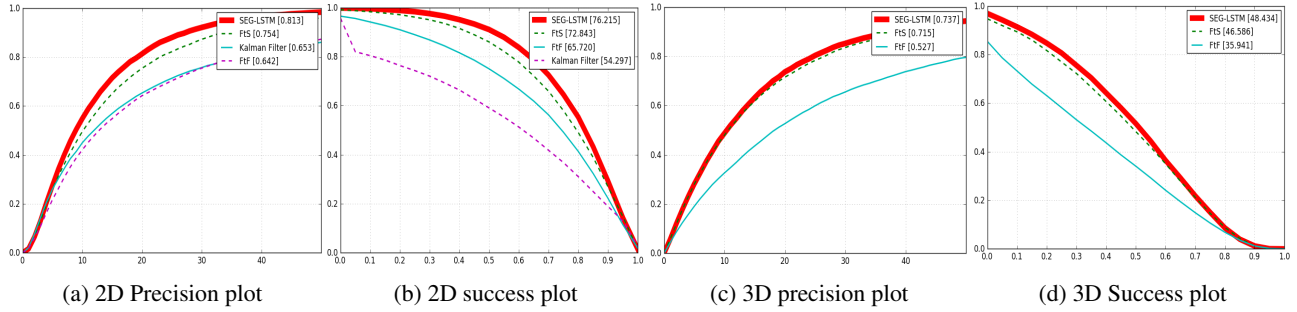


Figure 6: 2D and 3D space precision plot and success plot. x-axis for (a) to (d) are: pixel, IOU, centimeter and IOU.



Figure 7: Visualization of prediction results. Top two rows show the advantage of using semantics in temporal model: Red: ground truth; Blue: SEG-LSTM; Green: FtS; Yellow: FtF. SEG-LSTM can help detecting neighboring vehicles and restrict the viable path outside the space occupied by other vehicles (row 1); and confine the plausible car trajectory along the space above road (row 2). Row 3 and 4 show the benefit of incorporating social factors. Row 3 is the visualization of vanilla LSTM: the bottom left car predictions (denoted by red dots) are considered in isolation, moving dangerously close to the two front cars. Row 4 is the visualization of Social-LSTM: by considering two front cars in the scene, the prediction keeps the red dot car in a safer distance from the two front cars.

	Ave. disp. error	Final disp. error
Vanilla LSTM	0.0514	0.0880
<b>Social-LSTM</b>	<b>0.0494</b>	<b>0.0837</b>

Table 4: Results comparing vanilla LSTM and Social-LSTM in unit space. *Ave. disp. error* is the mean square error over all estimated points of a trajectory and the true points. *Final disp. error* is the distance between the predicted final destination and the true final destination at the end of the prediction period  $T_F$ .

other two temporal models. Fig. 7 shows the advantage of using SEG-LSTM: by utilizing the scene semantics (e.g., road, cars), it’s able to direct the recurrent net to output the most plausible path for vehicles and avoid collision with neighboring vehicles.

**Learning with social tensor:** On average, there is around 3 cars of interests in a high-way driving scene. Tab. 4 shows

that considering social interactions between cars improves the prediction accuracy over the isolation case. And we believe that in a more complex social scene such as urban driving, the improvement will be more pronounced.

## 5. Conclusions

In this paper, mimicking the human driver reaction time, we advocate to solve the problem of predicting traffic participants future trajectory given historical information. We propose a two-staged system incorporating various visual inputs with a temporally recurrent neural network. The recurrent network takes semantic information as an auxiliary input and further improve the predicting accuracy with the incorporation of social factor. Future works include incorporating vehicle pose from the RGB/depth input as extra source of “auxiliary information” and conducting experiments on the real-world images.



## References

- [1] Udacity. public driving dataset. <https://www.udacity.com/self-driving-car>, 2017.
- [2] A. Alahi, K. Goel, V. Ramanathan, A. Robicquet, L. Fei-Fei, and S. Savarese. Social lstm: Human trajectory prediction in crowded spaces. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- [3] A. Behl, O. H. Jafari, S. K. Mustikoveela, H. A. Alhaija, C. Rother, and A. Geiger. Bounding boxes, segmentations and object coordinates: How important is recognition for 3d scene flow estimation in autonomous driving scenarios? In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [4] S. Bengio, O. Vinyals, N. Jaitly, and N. Shazeer. Scheduled sampling for sequence prediction with recurrent neural networks. In *Advances in Neural Information Processing Systems*, 2015.
- [5] R. Chalasani and J. C. Principe. Deep predictive coding networks. *arXiv preprint arXiv:1301.3541*, 2013.
- [6] C. Chen, A. Seff, A. Kornhauser, and J. Xiao. Deepdriving: Learning affordance for direct perception in autonomous driving. In *Proceedings of the IEEE International Conference on Computer Vision*, 2015.
- [7] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*, 2014.
- [8] A. Clark. Whatever next? predictive brains, situated agents, and the future of cognitive science. *Behavioral and Brain Sciences*, 2013.
- [9] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- [10] M. Danelljan, G. Bhat, F. S. Khan, and M. Felsberg. Eco: Efficient convolution operators for tracking. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [11] M. Danelljan, G. Häger, F. S. Khan, and M. Felsberg. Discriminative scale space tracking. *IEEE transactions on pattern analysis and machine intelligence*, 2017.
- [12] M. Danelljan, A. Robinson, F. S. Khan, and M. Felsberg. Beyond correlation filters: Learning continuous convolution operators for visual tracking. In *European Conference on Computer Vision*. Springer, 2016.
- [13] M. Danelljan, K. F. S. Robinson, Andreas, and M. Felsberg. Beyond correlation filters: Learning continuous convolution operators for visual tracking. *European Conference on Computer Vision*, 2016.
- [14] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun. CARLA: An open urban driving simulator. In *Proceedings of the 1st Annual Conference on Robot Learning*, 2017.
- [15] T. Egner, J. M. Monti, and C. Summerfield. Expectation and surprise determine neural population responses in the ventral visual stream. *Journal of Neuroscience*, 2010.
- [16] M. Everingham, S. A. Eslami, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman. The pascal visual object classes challenge: A retrospective. *International journal of computer vision*, 2015.
- [17] C. Finn, I. Goodfellow, and S. Levine. Unsupervised learning for physical interaction through video prediction. In *Advances in Neural Information Processing Systems*, 2016.
- [18] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun. Vision meets robotics: The kitti dataset. *International Journal of Robotics Research*, 2013.
- [19] R. Girshick. Fast R-CNN. In *Proceedings of the IEEE International Conference on Computer Vision*, 2015.
- [20] R. B. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014.
- [21] A. Graves. Generating sequences with recurrent neural networks. *arXiv preprint arXiv:1308.0850*, 2013.
- [22] R. A. Hadi, G. Sulong, and L. E. George. Vehicle detection and tracking techniques: a concise review. *Signal & Image Processing*, 2014.
- [23] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask R-CNN. In *Proceedings of the IEEE International Conference on Computer Vision*, 2017.
- [24] K. He, X. Zhang, S. Ren, and J. Sun. Spatial pyramid pooling in deep convolutional networks for visual recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2015.
- [25] D. Held, S. Thrun, and S. Savarese. Learning to track at 100 fps with deep regression networks. In *European Conference on Computer Vision*. Springer, 2016.
- [26] J. F. Henriques, R. Caseiro, P. Martins, and J. Batista. High-speed tracking with kernelized correlation filters. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2015.
- [27] D. Hernandez-Juarez, L. Schneider, A. Espinosa, D. Vazquez, A. M. Lopez, U. Franke, M. Pollefeys, and J. C. Moure. Slanted stixels: Representing san francisco's steepest streets. In *British Machine Vision Conference*, 2017.
- [28] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 1997.
- [29] S. Hong, T. You, S. Kwak, and B. Han. Online tracking by learning discriminative saliency map with convolutional neural network. *International Conference on Machine Learning*, 2015.
- [30] J. Huang, V. Rathod, C. Sun, M. Zhu, A. Korattikara, A. Fathi, I. Fischer, Z. Wojna, Y. Song, S. Guadarrama, et al. Speed/accuracy trade-offs for modern convolutional object detectors. *arXiv preprint:1611.10012*, 2017.
- [31] A. Jain, A. R. Zamir, S. Savarese, and A. Saxena. Structural-rnn: Deep learning on spatio-temporal graphs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- [32] J. Kim and J. Canny. Interpretable learning for self-driving cars by visualizing causal attention. *International Conference on Computer Vision*, 2017.

972	[33]	W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y.	[50]	F. Yu, V. Koltun, and T. Funkhouser. Dilated residual net-	1026
973		Fu, and A. C. Berg. SSD: Single shot multibox detector. In		<i>Proceedings of the IEEE Conference on Computer</i>	1027
974		<i>European conference on computer vision</i> , 2016.		<i>Vision and Pattern Recognition</i> , 2017.	1028
975	[34]	C. Ma, J.-B. Huang, X. Yang, and M.-H. Yang. Hierarchical			1029
976		convolutional features for visual tracking. In <i>Proceedings</i>			1030
977		<i>of the IEEE International Conference on Computer Vision</i> ,			1031
978		2015.			1032
979	[35]	D. V. McGehee, E. N. Mazzae, and G. S. Baldwin. Driver			1033
980		reaction time in crash avoidance research: Validation of a			1034
981		driving simulator study on a test track. In <i>Proceedings of</i>			1035
982		<i>the human factors and ergonomics society annual meeting</i> ,			1036
983		2000.			1037
984	[36]	M. Mueller, V. Casser, J. Lahoud, N. Smith, and B. Ghanem.			1038
985		UE4Sim: A Photo-Realistic Simulator for Computer Vision			1039
986		Applications. <i>ArXiv e-prints</i> , Aug. 2017.			1040
987	[37]	M. Mueller, N. Smith, and B. Ghanem. A benchmark and			1041
988		simulator for uav tracking. In <i>European Conference on Com-</i>			1042
989		<i>puter Vision</i> , 2016.			1043
990	[38]	Y. Qi, S. Zhang, L. Qin, H. Yao, Q. Huang, and J. L. M.-H.			1044
991		Yang. Hedged deep tracking. In <i>Proceedings of the IEEE</i>			1045
992		<i>International Conference on Computer Vision</i> , 2016.			1046
993	[39]	J. Redmon, S. K. Divvala, R. B. Girshick, and A. Farhadi.			1047
994		You only look once: Unified, real-time object detection. In			1048
995		<i>Proceedings of the IEEE International Conference on Com-</i>			1049
996		<i>puter Vision</i> , 2016.			1050
997	[40]	S. Ren, K. He, R. Girshick, and J. Sun. Faster R-CNN: To-			1051
998		wards real-time object detection with region proposal net-			1052
999		works. In <i>Advances in neural information processing sys-</i>			1053
1000		<i>tems</i> , 2015.			1054
1001	[41]	G. Ros, L. Sellart, J. Materzynska, D. Vazquez, and A. M.			1055
1002		Lopez. The synthia dataset: A large collection of synthetic			1056
1003		images for semantic segmentation of urban scenes. In <i>Pro-</i>			1057
1004		<i>ceedings of the IEEE Conference on Computer Vision and</i>			1058
1005		<i>Pattern Recognition</i> , 2016.			1059
1006	[42]	E. Santana and G. Hotz. Learning a driving simulator. <i>arXiv</i>			1060
1007		<i>preprint arXiv:1608.01230</i> , 2016.			1061
1008	[43]	S. Thrun. Artificial intelligence for robotics course. <i>Udacity,</i>			1062
1009		<i>lecture note</i> , 2016.			1063
1010	[44]	T. Tieleman and G. Hinton. Lecture 6.5-rmsprop: Divide			1064
1011		the gradient by a running average of its recent magnitude.			1065
1012		<i>COURSERA: Neural networks for machine learning</i> , 2012.			1066
1013	[45]	A. Vemula, K. Muelling, and J. Oh. Social Attention: Mod-			1067
1014		eling Attention in Human Crowds. <i>International Conference</i>			1068
1015		<i>on Robotics and Automation</i> , 2018.			1069
1016	[46]	D. Wu, W. Zou, X. Li, and Y. Zhao. Kernalised multi-			1070
1017		resolution convnet for visual tracking. In <i>Proceedings of the</i>			1071
1018		<i>IEEE Conference on Computer Vision and Pattern Recogni-</i>			1072
1019		<i>tion workshop</i> , 2017.			1073
1020	[47]	Y. Wu, J. Lim, and M.-H. Yang. Online object tracking: A			1074
1021		benchmark. In <i>Proceedings of the IEEE International Con-</i>			1075
1022		<i>ference on Computer Vision</i> , 2013.			1076
1023	[48]	Y. Wu, J. Lim, and M.-H. Yang. Object tracking benchmark.			1077
1024		<i>IEEE Transactions on Pattern Analysis and Machine Intelli-</i>			1078
1025		<i>gence</i> , 2015.			1079
	[49]	H. Xu, Y. Gao, F. Yu, and T. Darrell. End-to-end learning of			
		driving models from large-scale video datasets. <i>IEEE con-</i>			
		<i>ference on computer vision and pattern recognition</i> , 2017.			

# Appendices

## A. Traffic participants detection

		conf	0.50	0.70	0.95
Faster-RCNN	precision		81.2	81.1	80.0
	recall		91.4	91.3	90.2
	f-score		86.0	86.0	84.8
SSD (NMS:0.60)	precision		82.1	81.6	79.2
	recall		<b>92.7</b>	92.0	89.4
	f-score		87.1	86.5	84.0
SSD (NMS:0.45)	precision		<b>92.0</b>	91.1	87.4
	recall		92.3	91.4	87.7
	f-score		<b>92.2</b>	91.3	87.6

Table 5: Comparison of SSD and Faster-RCNN for vehicle detection on the collected SYNTHIA dataset. It can be seen that: (1) both meta-architectures are robust to the cut-off confident thresholds. (2) SSD is more competitive in the collected dataset when targets of interest are comparatively plarge.

## B. Instance association via tracking

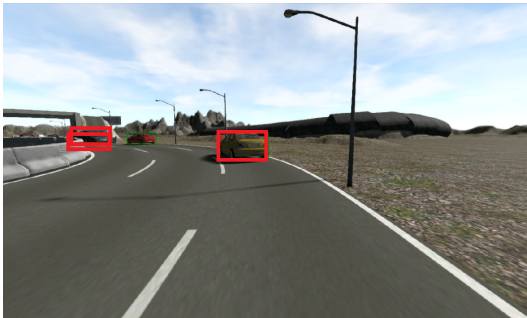


Figure 8: Instance association can help to disambiguate the false positive by enforcing temporal consistency requirement. Due to the high recall requirement for vehicle detection, there could be correspondingly higher false positive rate (*e.g.*, double detection of red boxed in the far left). Instance association across multiple frames will discard the spurious detection while remain temporal consistency.

## C. 2D and 3D occupancy grid

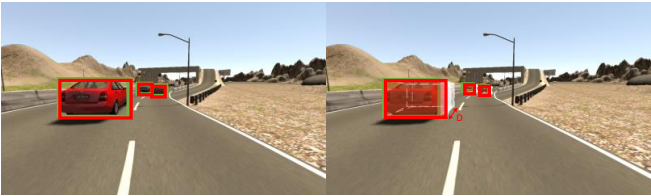


Figure 9: The evaluation of 2D and 3D prediction. 2D: The bounding box is set to cover the whole vehicle. 3D: The depth channel is combined with 2D bounding box to predict 3D occupancy.