

# 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 car’s future odometry 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 using privileged information (e.g., scene semantics) 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 [3]: 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 not just where other cars are, as in the localization case,

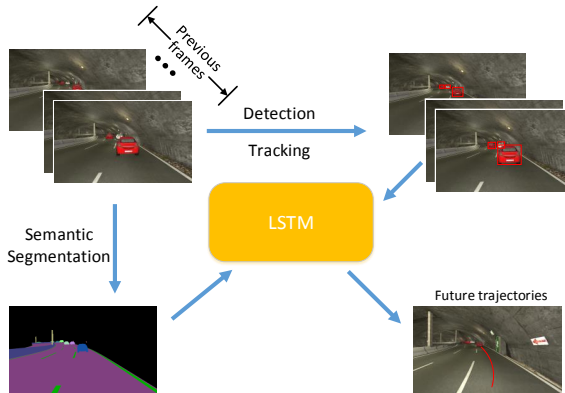


Figure 1: We propose a system to predict future trajectories of vehicles. Given historical information, obtained by detection, tracking and semantic segmentation, a multi-stream Long Short Term Memory (LSTM) based network is trained to generate future trajectory.

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 if stopping engaging the brake (remove foot from accelerator and apply brake). Accident reconstruction specialists commonly use 1.5 seconds [24]. Therefore, the ability to predict traffic participants’ future trajectories will greatly benefits 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, 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 driver makes decision by sequence of input frames instead of simply 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.*, [35, 36, 25]) 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 target’s future trajectory is simply a moot problem. However, traffic participants generally exhibit regular trajectories (*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 traditional visual tracking community and current autonomous driving research community, there is a lack of proper metrics for evaluating the temporal prediction result.

In our approach, there are three key components: traffic participants detection, instance tracking and future trajectory prediction as shown in Fig. 1. Our contributions are as follows:

- We verify the feasibility of predicting vehicles’ future trajectories for assisting mediated perception. The proposed system stems from the issue of human driver reaction time in the autonomous driving context.
- We present a novel “tracking-by-detection” framework for robust vehicle tracking. By updating tracker’s target representation via detection association, we also solve the problem of instance association between frames.
- We design various temporal models for the problem of predicting future trajectories based on historical data. Results show that the temporal model generating intermediate representation performs better than the frame-to-frame based temporal model.
- We demonstrate the effectiveness of using “privileged information” (*e.g.*, scene semantics) further boost the prediction accuracy.
- We construct a time-series dataset for vehicle trajectory prediction based on the SYNTHIA dataset [29] and formalize the problem into a 3D occupancy grid problem given depth information.

## 2. Related Work

### 2.1. Object detection

Vehicle detection is one key component of an autonomous driving system. Typical algorithms output bounding boxes on detected vehicles. 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 based on these networks – such as the line of works on the R-CNN [13], Fast R-CNN [12], Faster R-CNN [28], Mask R-CNN [14] and SSD [22] – are now good enough to be deployed in a consumer products. R-CNN [13] combines Selective Search and deep learning features to detect objects. SPPNet [15] exploits a spatial pyramid pooling layer to extract multi-scale deep features. Later Fast R-CNN [12] introduces the multi-task learning to fine-tune all layers in their network. At last, Faster R-CNN [28] proposes a region proposal network (RPN) and integrates RPN and Fast R-CNN to build an end-to-end architecture. Meanwhile, SSD [22] uses a series of boxes and directly gives the probability for each object class. In addition, YOLO [27] proposes a fast object detection method, which divides the input image into grids, then predicts object positions and confidences for each category in each grid.

### 2.2. Instance association

Instance association between frames is a prerequisite to solve the proposed intrinsically temporal problems. 3D scene flow is estimated in [2] by exploiting recognition to overcome the challenging scenarios in the present of large displacement or local ambiguities. However, a total runtime of 10 minutes per frame is way too prohibitive for the deployment in the context of autonomous driving. To accomplish the task of instance association between frames, we resort to the well studied field of visual tracking. Discriminative Correlation Filters (DCF) [17] have demonstrated excellent performance for high-speed generic visual object tracking. Built upon their seminal work, there has been a plethora of recent improvements [33, 32, 18, 23, 7, 16, 34, 5] relying on convolutional neural network (CNN) pre-trained on ImageNet as a feature extractor for visual tracking.

### 2.3. Recurrent neural networks

Recent advances in recurrent neural network for modeling visual sequential data are also related to our work. Mostly, RNN is trained as a pilot network in the context of behaviour reflex approaches. Long Short Term Memory (LSTM) cells are usually exploited as one intermediate component of an end-to-end driving model. A combination of a fully-convolutional network and an LSTM is proposed in [37] to learn from large-scale vehicle action data. Their goal is to predict future egomotion including multi-modal discrete and continuous driving behaviors. Inverse turning radius is predicted by a LSTM network in [20]: they generate a heat map of attention at each time step conditioned on the previous hidden states and a current convolutional feature cube to produce more succinct visual explanations. [9] used 3D convolutional layers to extract visual features, then fed them into LSTM layers to capture the sequential rela-

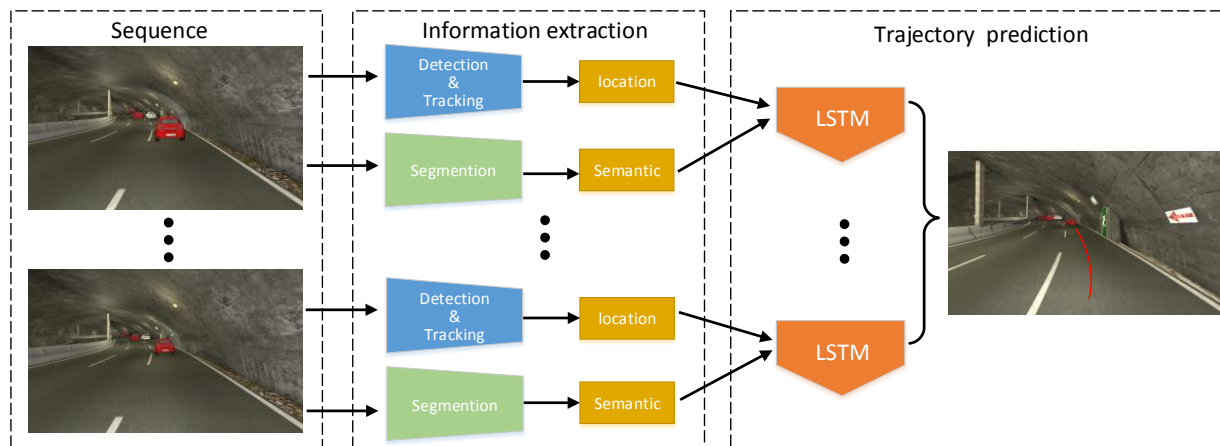


Figure 2: The framework for the proposed vehicle future trajectory prediction system. The input is the sequential data obtained from the host car. Historical information is extracted by detection & tracking. Along with scene semantics as an auxiliary information, multi-stream inputs are fed into our LSTM model to predict future trajectories.

tion. In addition to using deep features, [21] trained a large recurrent neural network using a reinforcement learning approach to map images directly to steering angles, with the purpose to keep the car on track.

## 2.4. Dataset for autonomous driving

Vision-based semantic segmentation in urban scene is a key functionality for autonomous driving. KITTI benchmark suite [11] 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 [4] 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 [30] 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. A synthetic dataset named SYNTHIA [25] is originally provided for urban scene semantic segmentation generation. Synthetic dataset has an obvious advantage: a fine annotated image in Cityscape 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 describe the overall framework for vehicle trajectory prediction and the details of implementation are presented in Sec 4.2. Our goal is via collecting visual information from the past frames to predict the traffic participants future trajectory. Fig. 2 shows the overall architecture of the proposed framework.

### 3.1. 3D Traffic Participants Trajectory Prediction

We propose to learn a generic approach from history information and formulate the problem as predicting future traffic participants' trajectories. Our work is most related to [37] where the problem was to predict future (next 1/3rd second) feasible actions of egomotion from a motion reflex perspective. We, instead, predict other traffic participants' future trajectories in a longer time span (next 1.6 second). Formally, our problem can be defined as a mapping  $\mathcal{M}$  between historical trajectory  $H = \{h_t^p, t = 1, \dots, T\}$  with auxiliary information  $A = \{a_t^p, t = 1, \dots, T\}$  for traffic participant  $p$  and the future trajectory  $J = \{j_t^p, t = T + 1, \dots\}$ :

$$J \leftarrow \mathcal{M}_p(h, a) : H \times A \quad (1)$$

where the traffic participants  $p$  can be of vehicles, cyclists, pedestrians, etc. In this paper, we focus on vehicle as the sole traffic participant category and ignore the suffix  $p$  from now on. Historical and future trajectories are defined in a 3D occupancy grid:

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

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

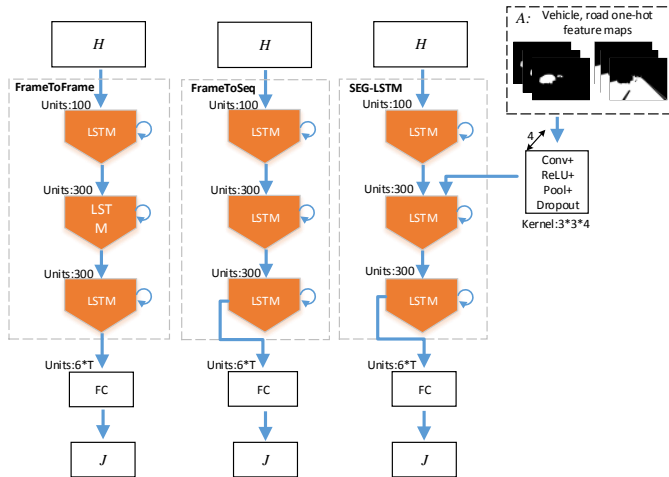


Figure 3: Illustrations of three temporal networks.  $H$  represents historical trajectories,  $J$  is the  $T$  time step trajectories to be predicted. Left: Frame-to-Frame (FtF); Middle: Frame-to-Sequence (FtS); Right: SEG-LSTM, with extra auxiliary information as input  $A$ .

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.

Auxiliary information  $A$  could be scene semantics to regularize the path of future trajectory or car pose information to decide whether it's an approaching car or a parallel driving car.

### 3.2. Sequence modeling

For predicting vehicle's future trajectory, we compare three temporal models based on LSTM cells 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 to better guide the trajectory prediction.

**Frame-to-Frame (FtF) model:** is akin to character prediction in a language model. 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  $j_t$  can be estimated as:

$$j_t \leftarrow \mathcal{M}_p(j_{t-1}, h) : H \quad (5)$$

**Frame-to-Sequence (FtS) model:** differs from the FtF model in that instead of sampling frame by frame for predicting future trajectories, FtS model directly predicts a trajectory sequence of length  $T$  (corresponding to the total driver reaction time). The assumption is 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  $j_{t,t+1,\dots,T}$  to be predicted can be estimated as:

$$j_{t,t+1,\dots,T} \leftarrow \mathcal{M}_p(h) : H \quad (6)$$

**SEG-LSTM model:** take advantages of scene semantics from historical frames as auxiliary information  $A$ . The second LSTM layer fuses the outputs from the bottom-most LSTM with the current frames's higher level semantic information. Semantic inputs are one-hot encoded images with channel number  $C$  as the number of semantic category. The higher level semantic is extracted from the semantic input through a four-layer CNN, producing a joint representation of the semantic auxiliary information and the historical trajectory. As described above, the 3D occupancy just has six dimensions:  $[x, y, w, h, d_{min}, d_{max}]$ . So the bottom-most LSTM not only can learn the temporal information between time-varying trajectory, but also can balance input dimensions of both trajectory information and the semantic information. Note that this model differs from that of [37] in that instead of using semantic segmentation as a side task for better directing the learning process of the convnets, SEG-LSTM directly takes scene semantics as input and directs the recurrent net future trajectory prediction. With auxiliary  $A$ , the trajectory  $j_{t,t+1,\dots,T}$  to be predicted can be estimated as:

$$j_{t,t+1,\dots,T} \leftarrow \mathcal{M}_p(h, a) : H \times A \quad (7)$$

## 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



## 1 Car trajectory ground truth collection

**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 step size 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.

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 car trajectory prediction dataset based on the SYNTHIA [29] 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 especially concern 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 [3], we set the reliable car perception as 30 meters so as to guarantee satisfactory control quality when the speed

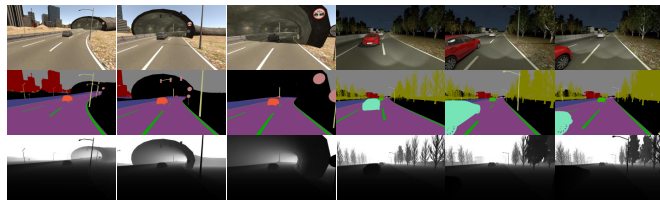


Figure 4: Examples of collected dataset for vehicle trajectory prediction. From top to bottom: RGB, semantic labeling, depth maps for three continuous frames under two dynamic light conditions.

of the host car does not exceed 72km/h. Since depth information is also provided alongside in the SYNTHIA dataset, we back-projecting 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 tracking vehicle. Hence, we start tracking the target when the  $POR$  of 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 [37] 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 [37] 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 [24].

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. We list our logic for collecting continuous car tracking in List 1. Some sample frames with semantic labels and depth information is shown in Fig. 4.

## 4.2. Implementation Details

### 4.2.1 Traffic participants detection

To reliably detect traffic participants, we compared two state-of-the-art approaches: SSD [22] and FasterRCNN [28]. We use the pretrained network on the Pascal VOC detection dataset [10] 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 [19] and we refer keen readers to the aforementioned paper for a more complete comparison to achieve the right speed/memery/accuracy balance for a given application

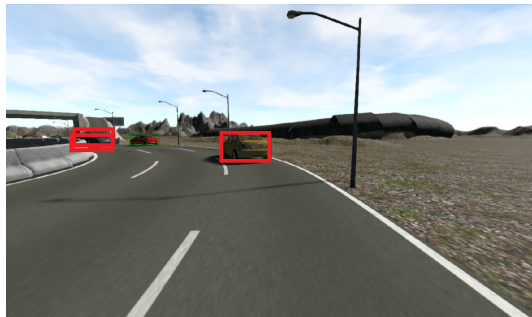


Figure 5: Instance association can help to disambiguate the false positive by enforcing temporal consistency requirement. Due to the high recall requirement for vehicle detection, there could be correspondingly higher false positive rate (e.g., double detection of red boxed in the far left). Instance association across multiple frames will discard the spurious detection.

and platform.

As it is also pointed out in the paper [19] that SSD, albeit is less advantageous in detecting small objects, it's very competitive for detecting larger objects. In this paper, we concerns mostly about cars that are close to the driver. Hence, given 0.1% POR as the the cut-off threshold, the presented problem favors detectors that are robust in detecting larger objects. Tab. 2 also verifies that SSD is indeed more competitive when objects of interest 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) threshold. It can be seen that both meta-architectures are robust to the cut-off confident scores. In terms of the effect of NMS cut-off threshold, allowing larger 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 confident score of 0.5 as the cut-off threshold and 0.45 as the default jaccard overlap for NMS.

#### 4.2.2 Instance tracking

To associate traffic participants across different frames (*i.e.*, instance tracking), we creatively combine detecting with tracking. The marriage between DCF [17], which has the advantage of being efficient in training translational images in the fourier space, and deep features, which excel at image representation, further advances the visual tracking community [26, 8, 34, 5]. However, in the pursuit of ever increasing tracking accuracy, their characteristic speed and real-time capability have gradually faded. In the context of autonomous driving, accurate scale estimation of a target is also a prerequisite. Most state-of-the-art methods employ

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

Table 1: Dataset statistics

an exhaustive scale search to estimate the target size which is computationally expensive and struggles when encountering with large scale variations. For this task, we adopt the real-time scale adaptive tracker fDSST [6] 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 this "tracking by detection" framework is essential for robust object tracking, especially when the objects undergoing rapid out of plane rotations. In return, the instance association scheme helps to alleviate the problem of false positive detection (*c.f.* Fig. 5). 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 only "car" and "road". The input semantics are one-hot encoded in the dimension of  $95 \times 160 \times 2$ . The CNN for learning high level semantics is composed of four 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. We adopt the state-of-the-art semantic segmentation networks: Dilated Residual Networks (DRN) [38] and retrain on the SYNTHIA corpus. DRN achieves 68% and 82% pixel mAP on the car and road categories.

#### 4.3. Evaluation

**Baseline: Kalman filter**, also known as linear quadratic estimation, is a popular technique for estimating the state of a system [31]. For the task of trajectory prediction, Kalman filters estimates a continuous state and gives a uni-modal distribution. The Kalman filter represents all distributions

	conf	0.5	0.55	0.6	0.65	0.7	0.75	0.8	0.85	0.9	0.95
Faster-RCNN	precision	81.2	81.2	81.1	81.1	81.1	81.0	81.0	80.7	80.5	80.0
	recall	91.4	91.4	91.3	91.3	91.3	91.3	91.3	90.9	90.7	90.2
	f-score	86.0	86.0	86.0	86.0	86.0	86.0	86.0	85.5	85.3	84.8
SSD (NMS:0.60)	precision	82.1	81.9	81.8	81.8	81.6	81.1	80.9	80.4	80.0	79.2
	recall	<b>92.7</b>	92.5	92.3	92.3	92.0	91.5	91.3	90.7	90.2	89.4
	f-score	87.1	86.9	86.7	86.7	86.5	86.0	85.8	85.2	84.8	84.0
SSD (NMS:0.45)	precision	<b>92.0</b>	91.8	91.7	91.4	91.1	91.0	90.5	90.2	89.4	87.4
	recall	92.3	92.1	92.0	91.8	91.4	91.3	90.8	90.5	89.7	87.7
	f-score	<b>92.2</b>	92.0	91.8	91.6	91.3	91.1	91.1	90.4	89.6	87.6

Table 2: 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 large.

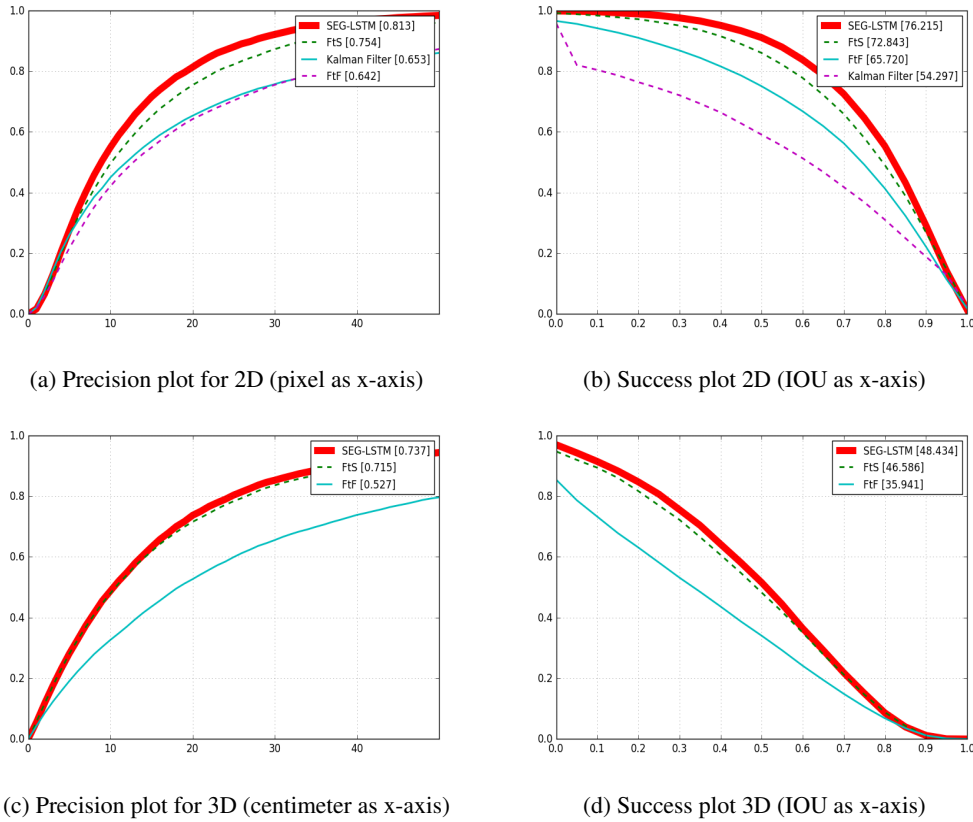


Figure 6: 2D and 3D space precision plot and success plot .

of the Gaussians and iterates two states: (1) measurement updates; (2) motion updates. In the baseline, we set the initial uncertainly covariance as  $1e4$ , measurement noise and motion noise are all set with unit value.

**2D Evaluation Methodology:** Following the evaluation strategy of [35], all trajectories are compared using two measures: precision and success. Precision is measured as the distance between the centers of the ground truth bound-

ing 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

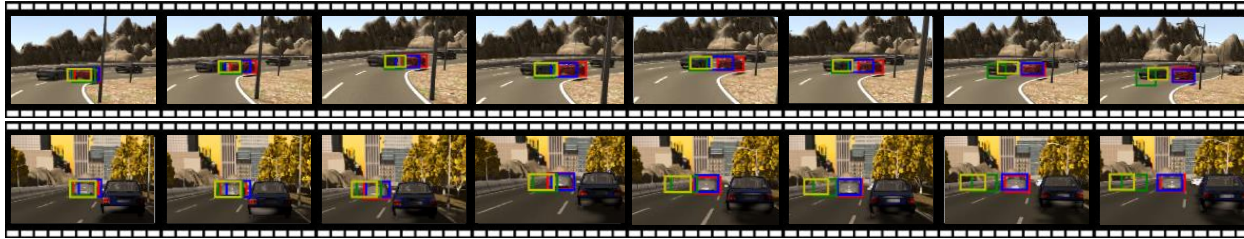


Figure 8: The prediction results in SYNTHIA dataset. Red: ground truth; Blue: SEG-LSTM; Green: FtS; Yellow: FtF. The results using SEG-LSTM network show that by utilizing the scene semantics (e.g., road, cars), it’s able to direct the predicted trajectory from the recurrent net within the most plausible path. Top: SEG-LSTM can help detecting neighbouring vehicles and restricts the viable path outside the space occupied by other vehicles. Bottom: SEG-LSTM confines the possible car trajectories along the space above road.

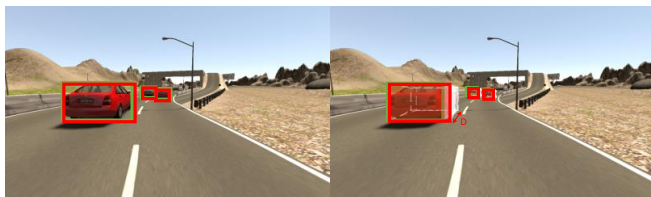


Figure 7: The evaluation of 2D and 3D prediction. 2D: It is the bounding box to cover the whole vehicle. 3D: The depth channel is combined with 2D bounding box to predict trajectory.

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. 3 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 coordinate. The auxiliary learning with semantic information in SEG-LSTM also consistently outperforms other two temporal models. Fig. 8 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.

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 3: Results of 2D performance evaluation 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 4: Results of 3D performance evaluation on mean coverage (higher is better) and center error (lower is better).

## 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. Future works include incorporating vehicle pose from the RGB/depth input as extra source of “privileged 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. Behl, O. H. Jafari, S. K. Mustikovela, 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.
- [3] 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.
- [4] 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.
- [5] 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.
- [6] M. Danelljan, G. Häger, F. S. Khan, and M. Felsberg. Discriminative scale space tracking. *IEEE transactions on pattern analysis and machine intelligence*, 2017.
- [7] 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.
- [8] 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.
- [9] S. Du, H. Guo, and A. Simpson. Self-driving car steering angle prediction based on image recognition. In *Stanford: CS231n*, 2017.
- [10] 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.
- [11] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun. Vision meets robotics: The kitti dataset. *International Journal of Robotics Research*, 2013.
- [12] R. Girshick. Fast R-CNN. In *Proceedings of the IEEE International Conference on Computer Vision*, 2015.
- [13] 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.
- [14] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask R-CNN. In *Proceedings of the IEEE International Conference on Computer Vision*, 2017.
- [15] 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.
- [16] 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.
- [17] 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.
- [18] 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.
- [19] 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.
- [20] J. Kim and J. Canny. Interpretable learning for self-driving cars by visualizing causal attention. *Proceedings of the IEEE International Conference on Computer Vision*, 2017.
- [21] J. Koutník, G. Cuccu, J. Schmidhuber, and F. Gomez. Evolving large-scale neural networks for vision-based torcs. 2013.
- [22] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg. SSD: Single shot multibox detector. In *European conference on computer vision*, 2016.
- [23] C. Ma, J.-B. Huang, X. Yang, and M.-H. Yang. Hierarchical convolutional features for visual tracking. In *Proceedings of the IEEE International Conference on Computer Vision*, 2015.
- [24] D. V. McGehee, E. N. Mazzae, and G. S. Baldwin. Driver reaction time in crash avoidance research: Validation of a driving simulator study on a test track. In *Proceedings of the human factors and ergonomics society annual meeting*, 2000.
- [25] M. Mueller, N. Smith, and B. Ghanem. A benchmark and simulator for uav tracking. In *European Conference on Computer Vision*, 2016.
- [26] Y. Qi, S. Zhang, L. Qin, H. Yao, Q. Huang, and J. L. M.-H. Yang. Hedged deep tracking. In *Proceedings of the IEEE International Conference on Computer Vision*, 2016.
- [27] J. Redmon, S. K. Divvala, R. B. Girshick, and A. Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE International Conference on Computer Vision*, 2016.
- [28] S. Ren, K. He, R. Girshick, and J. Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, 2015.
- [29] G. Ros, L. Sellart, J. Materzynska, D. Vazquez, and A. M. Lopez. The synthia dataset: A large collection of synthetic images for semantic segmentation of urban scenes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- [30] E. Santana and G. Hotz. Learning a driving simulator. *arXiv preprint arXiv:1608.01230*, 2016.
- [31] S. Thrun. Artificial intelligence for robotics course. *Udacity, lecture note*, 2016.
- [32] L. Wang, W. Ouyang, X. Wang, and H. Lu. Visual tracking with fully convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition workshop*, 2015.
- [33] L. Wang, W. Ouyang, X. Wang, and H. Lu. Stct: Sequentially training convolutional networks for visual tracking. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition workshop*, 2016.

972		1026
973		1027
974		1028
975		1029
976	[34] D. Wu, W. Zou, X. Li, and Y. Zhao. Kernalised multi-	1030
977	resolution convnet for visual tracking. In <i>Proceedings of the</i>	1031
978	<i>IEEE Conference on Computer Vision and Pattern Recognition</i>	1032
979	<i>workshop</i> , 2017.	1033
980	[35] Y. Wu, J. Lim, and M.-H. Yang. Online object tracking: A	1034
981	benchmark. In <i>Proceedings of the IEEE International Con-</i>	1035
982	<i>ference on Computer Vision</i> , 2013.	1036
983	[36] Y. Wu, J. Lim, and M.-H. Yang. Object tracking benchmark.	1037
984	<i>IEEE Transactions on Pattern Analysis and Machine Intelli-</i>	1038
985	<i>gence</i> , 2015.	1039
986	[37] H. Xu, Y. Gao, F. Yu, and T. Darrell. End-to-end learning of	1040
987	driving models from large-scale video datasets. In <i>Proceed-</i>	1041
988	<i>ings of the IEEE Conference on Computer Vision and Pattern</i>	1042
989	<i>Recognition</i> , 2017.	1043
990	[38] F. Yu, V. Koltun, and T. Funkhouser. Dilated residual net-	1044
991	works. In <i>Proceedings of the IEEE Conference on Computer</i>	1045
992	<i>Vision and Pattern Recognition</i> , 2017.	1046
993		1047
994		1048
995		1049
996		1050
997		1051
998		1052
999		1053
1000		1054
1001		1055
1002		1056
1003		1057
1004		1058
1005		1059
1006		1060
1007		1061
1008		1062
1009		1063
1010		1064
1011		1065
1012		1066
1013		1067
1014		1068
1015		1069
1016		1070
1017		1071
1018		1072
1019		1073
1020		1074
1021		1075
1022		1076
1023		1077
1024		1078
1025		1079