



HAL
open science

Applying map-masks to Trajectory Prediction for Interacting Traffic-Agents

Vyshakh Palli Thazha, David Filliat, Javier Ibañez-Guzmán

► **To cite this version:**

Vyshakh Palli Thazha, David Filliat, Javier Ibañez-Guzmán. Applying map-masks to Trajectory Prediction for Interacting Traffic-Agents. Deep Learning for Automated Driving (DLAD) workshop, IEEE International Conference on Intelligent Transportation Systems, Oct 2019, Auckland, New Zealand. hal-02393250

HAL Id: hal-02393250

<https://hal.science/hal-02393250>

Submitted on 4 Dec 2019

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Applying map-masks to Trajectory Prediction for Interacting Traffic-Agents

Vyshakh Palli-Thazha, David Filliat and Javier Ibañez-Guzmán

Abstract—Autonomous vehicles perceive their surroundings and foresee the future behaviour of all other relevant interacting traffic-agents in order to navigate safely and to operate in the public road networks. Trajectory prediction is difficult due to the stochastic manner these different types of traffic-agents interact with each other and the changing navigation context. In this work, we propose a recurrent artificial neural network LSTM encoder-decoder based architecture to predict the movement of traffic agents. It uses map-masks of the area surrounding the ego vehicle and previous trajectory information to predict the trajectory of interacting traffic agents. This paper compares the proposed approach with LSTM baselines, using the NuScenes dataset which includes LiDAR point-cloud ground-truth data for traffic agents plus map information. Experimental results show that the proposed method outperforms the baselines based on the prediction accuracy.

I. INTRODUCTION

Sequence prediction problems have been an important area of research with multiple applications like Natural Language Processing, Weather Forecast, Economics and Intelligent Systems. Autonomous vehicle systems come across this problem when trying to understand the behaviour of all other traffic-agents that it interacts with. This could be a pedestrian, car, bike or any object that appear in its navigation space. Predicting the trajectory of the involved agents allow the autonomous vehicle to prepare itself for future decisions and manoeuvres, after assessing the risk involved in the current scenario. The goal is to try and make the vehicle understand the environment like humans do. The high level of interactions, occlusions and changing context makes it hard to solve this prediction problem with high accuracy. Therefore multiple inputs, patterns and behaviours must be studied to come to an operational and acceptable accurate prediction.

Researchers have addressed the problem focusing on different aspects. Sequence predictions for single-agent behaviours have been solved using classical techniques but in real life this is unusually the case. There are multiple agents and their behaviour changes according to the class of the agent - whether it is a car, a pedestrian or a bike. Our research takes the class of the agent into account in the problem statement and tries to predict their different

behaviours. Another factor that plays an important role in real life traffic scenarios are the interactions between different classes of agents. Our current research also gives importance to this and has incorporated an interaction-aware model. Thirdly, our model takes into account, the surroundings of the ego-motion car, by using a map-mask that provides information on the driveable and non-driveable areas for different classes of agents. These map-masks also provide information on the path that could be taken by different agents.

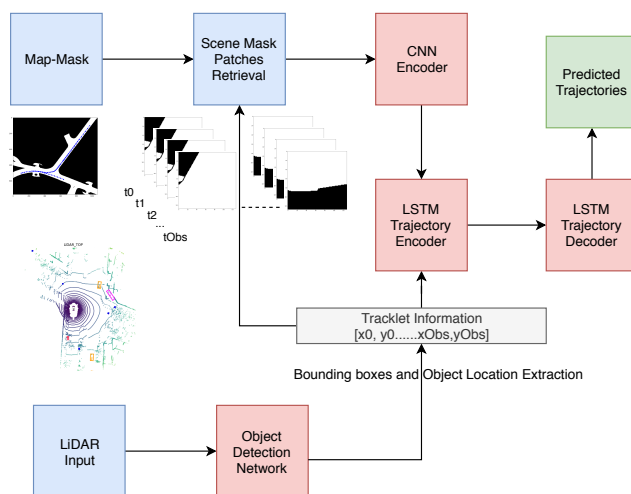


Fig. 1: System Architecture: The prediction module has two inputs: sequential historical trajectories of all traffic-agents and the associated map-mask patches of each agent for the observed time instances

In the recent years, the boom in high-computational-power processors and GPUs has lead to faster research and development cycles. More and more data-sets are being released for different applications - KITTI, Apolloscape, NuScenes and NGSIM. This has lead to a focus on deep learning algorithms which solve classical problems with extremely high rate of success and accuracy. This success of Learning Algorithms have attracted research on several methods like Convolutional Neural Networks(CNN), Recurrent Neural Networks(RNN), Long Short-Term Memory(LSTM) and Generative Adversarial Networks(GANs). RNNs and LSTMs have shown promising results with solving sequence prediction problems. We use one such LSTM encoder-decoder architecture (see Figure 3 for an overview) as the baseline for our study and use the NuScenes dataset [1] which provides 3D bounding-

Vyshakh Palli-Thazha is a PhD candidate at ENSTA Paris and Groupe Renault, Paris, France vyshakh.palli-thazha@ensta-paris.fr, vyshakh.v.palli-thazha@renault.com
David Filliat, Director, U2IS Lab, ENSTA Paris, Institut Polytechnique de Paris david.filliat@ensta-paris.fr
Javier Ibañez-Guzmán, Groupe Renault javier.ibanez-guzman@renault.com

boxes in LiDAR point-clouds for different traffic agents, object/agent IDs and do a preliminary study for improving trajectory prediction using map-masks (binary maps) for the environment which we use in our experiments.

In this work, the use of map-mask patches to improve the prediction of trajectories for different classes of interacting traffic-agents is proposed. Specifically, our contributions include:

- A new LSTM encoder-decoder architecture that uses Map-Mask patches to make trajectory predictions for different classes of traffic agents in drivable and non-drivable areas
- We evaluate the proposed model in comparison with LSTM baselines and is found to be superior in performance both in single-agent and multi-agent interaction scenarios.

The remainder of the paper is organised as follows: Section II mentions relevant research papers and results that we have taken inspiration from or has guided our research. Research related to single-agent predictions, multi-agent predictions, interaction-aware models, context-aware models are discussed in this section. Section III describes our map-mask based LSTM encoder-decoder architecture and the working of our complete trajectory prediction approach. Section IV discusses the results and the performance of our model and finally section V concludes the paper and mentions the future work.

II. RELATED WORKS

A. Classical Methods

Trajectory prediction for pedestrians, vehicles and other traffic agents has been studied using classical techniques such as Constant Velocity, Constant Acceleration, Linear Regressions, Kalman Filters [2], Monte Carlo Simulation [3], Time-series methods and Hidden Markov Models [4]. Most of these predictions are limited to short term predictions and even if these are long-term, their accuracy is affected because of the lack of context information and from ignoring the interaction between different classes of traffic-agents. A survey on motion prediction can be found in [5] where Lefèvre *et. al.* have analysed different trajectory prediction methods based on model completeness and risk assessment.

B. Learning based Models

Trajectory prediction methods have been applied to different sensor outputs, some solely with vision based detection outputs, others with 3D bounding-boxes obtained from LiDAR point-clouds. We base our work with the detections on 3D point-cloud data. Much research has gone into developing deep learning algorithms to infer such predictions. Alché *et. al.* [6] uses LSTMs to predict the future trajectory of a single target vehicle. It predicts the longitudinal and lateral trajectories of vehicles on a highway. In [7] A. Milan *et. al.* proposed an end-to-end learning approach for online multi-target tracking which uses RNNs to predict the state of each target step by step and uses LSTMs to achieve data association. They work on vision data obtained from

the 3D MOTChallenge dataset. Nikhil *et. al.* [8] uses a Convolutional Neural Network based approach to predict trajectories of surrounding vehicles. This model utilises highly parallelisable convolutional layers to handle temporal dependencies. The trajectory histories are embedded to a fixed size tensor and stacked convolutional layers are used to ensure temporal consistency. Jawed *et. al.* [9] also uses a similar architecture based on convolutional layers.

C. Interaction and/or Context Aware Models

There are several examples in the literature that tries to solve the prediction problem which involves interactions between traffic agents. Some of them focus on pedestrians alone or vehicles alone while others try to solve cases with interaction between multiple classes of agents.

Several models study interaction between vehicles: In [10] Hu *et. al.* uses a generative model to jointly predict the sequential motions of each pair of interacting vehicles and it uses as input past trajectories and environment information. Kim *et. al.* [11] uses occupancy grid and RNNs to predict the trajectory of vehicles on highways. In [12] Deo *et. al.* uses an LSTM encoder-decoder architecture that uses convolutional social pooling for learning inter-dependencies between vehicles. Here the spatial configuration of the agents in a scene is embedded into a spatial grid around the ego-motion vehicle and this grid is passed into convolutional and pooling layers to obtain the social context encoding. This along with an LSTM encoding of the agents' trajectory are concatenated and passed through a decoder to obtain the final predicted trajectory. A similar architecture is used by Park *et. al.* [13] with the decoder producing K likely trajectories over an occupancy grid using the beam search technique.

Some models study pedestrian interactions: Kooij *et. al.* [14] proposes a Dynamic Bayesian Network which uses the situational awareness and spatial layout to predict pedestrian paths. It uses pedestrian head orientation, distance between vehicles and pedestrians at expected point of closest approach and distance of pedestrian to curbside as latent states on a Switching Linear Dynamical System (SLDS) to predict changes in pedestrian dynamics. In [15] Alahi *et. al.* tries to learn general human movement and predicts their future trajectories. They model pedestrian motion by a pooling based LSTM architecture and jointly predicts the trajectories of all the people in a scene.

Models also use various context information to improve trajectory prediction: Lee *et. al.* [16] applies a Conditional Variational AutoEncoder (CVAE) based RNN encoder-decoder to make prediction for interacting agents. They use past trajectories as input along with a feature map, generated from scene elements like roads and sidewalks, that provides context information to the prediction model. Habibi *et. al.* [17] uses relative distance to curbside and status of pedestrian traffic lights as additional information to provide context to predict the pedestrian path. Ridel *et. al.* [18] implicitly models pedestrian interactions with vehicles to predict pedestrian behaviour. They use pedestrian head orientation as an input along with past trajectories. Another

interesting study is done by Ma *et. al.* [19] which uses 4D graphs to model interactions and classes of the interacting agents.

Our approach addresses the problem by using map-masks to facilitate the understanding of the scene environment. This study is a first step towards using semantic information in the form of maps for improving trajectory prediction. The trend is to add more information to maps. Works like OpenStreetMaps and NuScenes are releasing maps with information other than road geometry. Details such as cross-walk, sidewalk, traffic lights, stop lines, lanes are being added to maps. These semantics could be used to improve the perception of the environment.

III. PROBLEM FORMULATION AND METHOD

In a driving scenario, several objects might be present, sharing the same ego-vehicle work-space it navigates. These include pedestrians, other vehicles or cyclists. State-of-the-art perception algorithms provide bounding boxes/segmentation for each class of objects. There are algorithms which can perceive the said objects from camera images or from LiDAR point-cloud data. In our case, we use LiDAR point-cloud data and off-the-shelf detection algorithms to extract the features of all the traffic-agents within the field of view of the sensor used. For experimental purposes, the PointPillars object-detection network by Lang *et. al.* [20] is adopted. This network is chosen as it shows superior performances over other relevant frameworks and they have released their code on GitHub. The observed trajectory is taken from the output of the object detection framework. For training the prediction network, the ground truth trajectory available for the situation from the NuScenes dataset is used.

The feature set of each traffic agent k at time t can be defined by $f_t^k = (p_t^k, c_t^k, M_t^k)$ where $p_t^k = [x_t^k, y_t^k]$, x and y are the spatial coordinates of the detected traffic-agent and c is the class of the traffic agent, $c_i \in (1, 2, 3)$ where 1 is for pedestrians, 2 is for cars and 3 is for cyclists and M_t^k is a 128×128 pixels binary map patch of the area surrounding each detected object. The task is to observe the features of all traffic-agents in the interval $[1 : T_{obs}]$ and predict their positions in $[T_{obs}+1 : T_{pred}]$.

A. Trajectory Prediction using LSTMs

The prediction framework developed by T. Fernando *et. al.* [21] that applies LSTMs for predicting trajectories is applied, along with the training methods adopted by Bahdanau *et. al.* in their Neural Machine Translation work [22]. For each time-step, an object-pool is created, associated with a trajectory prediction module. It regresses short term and long-term trajectories for each traffic-agent in the object-pool. This is based on the historic trajectory of that agent and neighbouring trajectories. Therefore, each traffic-agent that resides in the pool is associated with a probable short term and long-term trajectory. This approach is advantageous because raw detection outputs can be noisy due to occlusions, false alarms, inaccurate bounding boxes, and missing detec-

tions. Studies show that occlusions last around 3-10 frames [23].

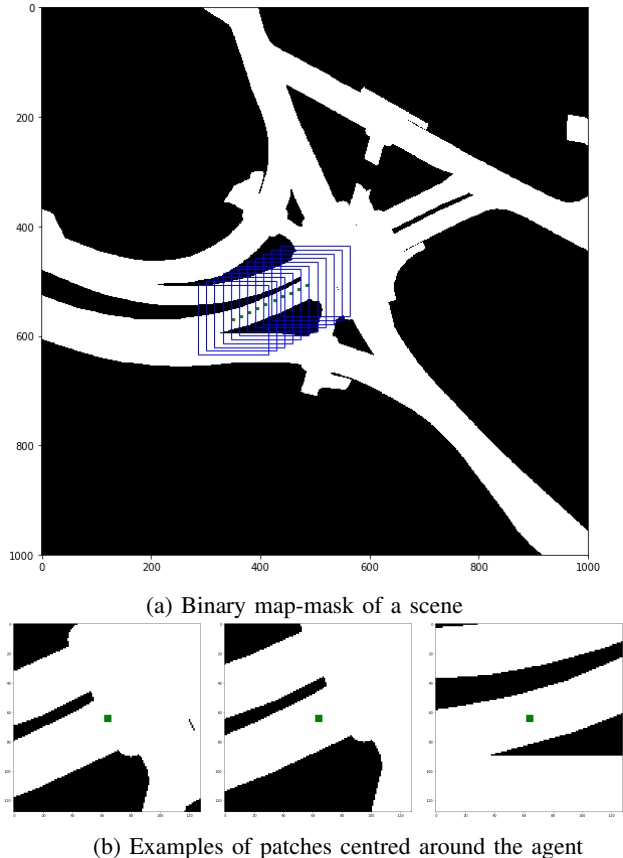


Fig. 2: Binary map-mask and map patches: A trajectory and the map-patches surrounding the agent are marked by blue squares, centred around the traffic agent marked in green

Trajectory prediction:

Let the historical trajectory of a traffic-agent k , from time-step 1 to time-step T_{obs} be given by,

$$p^k = [p_1, \dots, p_{T_{obs}}],$$

where,

$$p_i = [x_i, y_i]$$

for each time-step i . p_i composes of points in a Cartesian grid obtained from the top-down view of the LiDAR point-cloud detections. These historical trajectories are passed through the LSTM encoder of each respective traffic-agent to generate its historical embeddings as follows,

$$h_t^k = LSTM(p_t^k, h_{t-1}^k),$$

generating a sequence of historical embeddings.

A historical context vector $C_t^{H,k}$ is defined to encode the trajectory information from the traffic-agent of interest (k), which can be computed as a weighted sum of hidden states from $t = [1 : T_{obs}]$,

$$C_t^{H,k} = \sum_{j=1}^{T_{obs}} \alpha_{tj} h_j^k$$

and the weight α_{tj} as shown in [22] can be computed by,

$$\alpha_{tj} = \frac{\exp(e_{tj})}{\sum_{l=1}^T \exp(e_{tl})}$$

$$e_{tj} = a(h_{t-1}^k, h_j^k)$$

where the function a is a feed forward neural network which is trained jointly with the whole network.

The spatial context vector $C^{S,k}$ is used for embedding the neighbouring trajectories. The spatial weights, denoted by w_j^n , can be computed as,

$$w_j^n = \frac{1}{\text{dist}(n, j)}$$

where $\text{dist}(n, j)$ is the Euclidean distance between the n^{th} neighbour and the traffic-agent of interest at the j^{th} time-step, and w_j^n is the generated spatial attention weight. When there are N neighbouring trajectories in the local neighbourhood, and h_j^n is the encoded hidden state of the n^{th} neighbour at the j^{th} time-step, then the context vector for the spatial model is defined as,

$$C^{S,k} = \sum_{n=1}^N \sum_{j=1}^{T_{obs}} w_j^n h_j^n$$

The merged context vector, $C_t^{*,k}$, computed by,

$$C_t^{*,k} = \tanh([C^{H,k}; C^{S,k}])$$

is then passed through the LSTM decoder to predict the future trajectory for the traffic-agent of interest. This is explained in Section III-C.

It is by modifying the context vector that the information from the binary map-masks is input to the prediction model. This is explained in the next sub-section.

B. Map Input for Trajectory Prediction

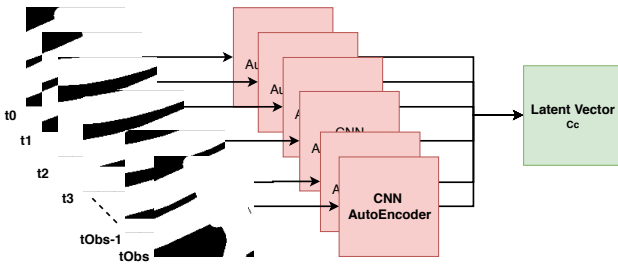


Fig. 3: Passing the map-patches through a CNN AutoEncoder generates the latent vector that stores information about the drivable area and this vector is passed on to the LSTM encoder architecture

For understanding the behaviour of different classes of traffic-agents in drivable and non-drivable spaces, the input of map information in the form of binary patches is studied. We take square patches of the size 128*128 pixels (see Figure 2) around each detected traffic-agent in the perception space

(see Figure 2a) of the ego-motion vehicle and for each time-step in the observed time interval. Each of these patches are passed through a Convolutional AutoEncoder producing a latent vector C_t^C which holds the information about the drivable space associated with the respective traffic-agent (see Figure 3).

This vector is then concatenated with the merged context vector $C_t^{*,k}$ from the LSTM Model to obtain the new latent vector $C_t^{M,k}$.

C. LSTM Decoder

Now, the final predicted trajectory q_t is obtained by passing the context vector $C_t^{M,k}$ in the case of map input or $C_t^{*,k}$ in the case of only LSTMs, through the LSTM decoder and two fully connected layers. The overall system architecture is illustrated in Figure 1.

$$q_t = LSTM_c(h_{t-1}^k, q_{t-1}, C_t^{M/*,k})$$

q_t^k is composed of points in a Cartesian grid. The decoder used for each traffic-agent class is different. Hence, based on c , we choose the LSTM to be used to decode. We let $T_{obs} = 3$ for short term predictions and $T_{obs} = 10$ for long term trajectory predictions.

$$q_t = [x_{Obs+1}, y_{Obs+1}, \dots, x_{Obs+N}, y_{Obs+N}]$$

where N is the prediction horizon.

The CNN AutoEncoder is tested both with and without pre-training and no significant difference has been found. Therefore, this model is attached to the LSTM encoder-decoder and trained end-to-end with Adam optimiser and a learning rate of $1 \times e^{-4}$ for 100 epochs and fine-tuned with a learning rate of $1 \times e^{-5}$.

IV. RESULTS AND EVALUATION

The prediction accuracy is measured based on two criteria as described in [24], [25]: Average Displacement Error (ADE) and Final Displacement Error (FDE) in pixels with respect to each time step t within the prediction horizon:

- 1) Average Displacement Error (ADE): Average distance between ground truth and our prediction over all predicted time-steps for the traffic agent.

$$ADE(i) = \frac{1}{T} \sum_{j=1,2,\dots,T} \sqrt{(\hat{x}_i^j - x_i^j)^2 + (\hat{y}_i^j - y_i^j)^2}$$

$$ADE = \frac{1}{n} \sum_{i=1,2,\dots,n} ADE(i)$$

- 2) Final Displacement Error (FDE): The distance between the predicted final destination and the true final position of the agent.

$$FDE(i) = \frac{1}{n} \sqrt{(\hat{x}_i^T - x_i^T)^2 + (\hat{y}_i^T - y_i^T)^2}$$

$$FDE = \frac{1}{n} \sum_{i=1,2,\dots,n} FDE(i)$$

where n is the total number of interacting agents in the test set, x_{ij} and y_{ij} denote the coordinates of the i^{th} agent in

the predicted time-step j and T denotes the final predicted time-step.

As training data, the NuScenes dataset is used. From this we obtain 3200 different trajectories with enough length (at least 30 time frames) for the experiments. These trajectories are obtained from 850 scenes of 20s recordings from the NuScenes data-collection car setup. These are split into train, validation and test sets of size 2100, 550 and 550 trajectories respectively. For single-agent trajectory prediction, these sequences are considered separately, and for the multi-agent trajectory predictions, the information of all interacting agents are taken into account. Each trajectory sequence is associated with the relevant interacting agent information also.

In order to illustrate the interest of using drivable/non-drivable map information, four variants of our approach are implemented.

- **Method 1:** Single Agent LSTM Model for all traffic agents: Separate LSTM models are used to predict the motion of single traffic agents. No interaction is involved.
- **Method 2:** Single Agent with LSTM Model and Map-Mask input: Map-mask is taken as input to predict the trajectories of single agents, not considering interactions.
- **Method 3:** Multi-Agent interactions and Model with LSTM only: Interactions are taken into consideration and only LSTM encoder-decoder architecture is used. This is similar to the method in [26].
- **Method 4:** Multi-Agent interactions and Model with LSTM and Map-Mask input: Multi-agent interactions are taken into account in the LSTM encoder-decoder architecture and map-masks are used to improve the predictions of traffic-agents involved in such interactions.

| Method | ADE(Meters) | FDE(Meters) |
|----------------------------------|-------------|-------------|
| 1. Single Agent LSTM only | 5.31 | 6.2 |
| 2. Single Agent LSTM + Map-Input | 1.67 | 1.95 |
| 3. Multi-Agent LSTM only | 4.34 | 5.21 |
| 4. Multi-Agent LSTM + Map-Input | 1.45 | 2.2 |

TABLE I: Comparison of Prediction Accuracy for different models

The performance of the model in different scenarios is shown in Figure 4. Figure 4a shows an example of straight trajectories and the prediction giving results in the non-drivable area when only LSTMs are used and better predictions confined to the drivable area when map-mask is added as input. Also the fitting of the trajectory prediction to road curves is better predicted when map-masks are introduced as can be seen in 4b. Another example of prediction being made in the drivable area is seen in Figure 4c.

The comparison of the ADE/FDE values for the LSTM baselines and the models with map-mask input is presented in Table I.

As can be seen in each case, single agent and multi-agent predictions, the addition of the map-mask patches

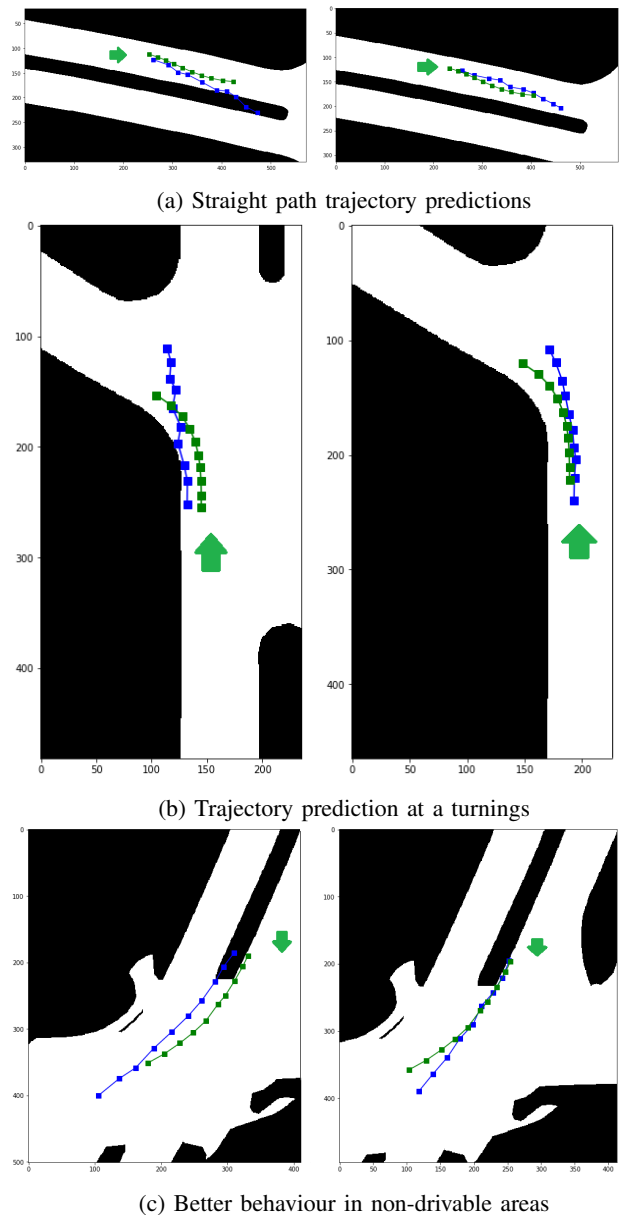


Fig. 4: Behaviour of Trajectory Prediction Models without(Left) and with Map-mask(Right) input. The green arrows in the images show the direction of prediction. The green path is the Ground truth and the blue path is the predicted trajectory

help the system make more accurate predictions. This gives us confidence to study the use of a map with much more information in the future to improve the predictions even more.

V. CONCLUSION AND FUTURE WORK

The proposed approach discusses and experimentally verifies the idea of using binary map-masks to improve the prediction of trajectories for different traffic agents. It compares single-agent and multi-agent trajectory prediction models and asserts the overall improvement of the prediction when map-masks are given as input. This study serves as a preliminary

step in the direction of using more environmental context information for improving prediction accuracy. As a future work, we plan to study OpenStreetMaps and use the multitude of information that can be obtained from their wide network of mapped areas.

REFERENCES

- [1] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liong, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom, "nusenes: A multimodal dataset for autonomous driving," *arXiv preprint arXiv:1903.11027*, 2019.
- [2] R. E. Kalman, "A new approach to linear filtering and prediction problems," *ASME Journal of Basic Engineering*, 1960.
- [3] S. Danielsson, L. Petersson, and A. Eidehall, "Monte carlo based threat assessment: Analysis and improvements," 07 2007, pp. 233 – 238.
- [4] J. Firl, H. Stubing, S. Huss, and C. Stiller, "Predictive maneuver evaluation for enhancement of car-to-x mobility data," 06 2012, pp. 558–564.
- [5] S. Lefèvre, D. Vasquez, and C. Laugier, "A survey on motion prediction and risk assessment for intelligent vehicles," *ROBOMECH Journal*, vol. 1, no. 1, p. 1, 2014. [Online]. Available: <https://hal.inria.fr/hal-01053736>
- [6] F. Alché and A. de La Fortelle, "An LSTM network for highway trajectory prediction," *CoRR*, vol. abs/1801.07962, 2018. [Online]. Available: <http://arxiv.org/abs/1801.07962>
- [7] A. Milan, S. H. Rezatofghi, A. R. Dick, K. Schindler, and I. D. Reid, "Online multi-target tracking using recurrent neural networks," *CoRR*, vol. abs/1604.03635, 2016. [Online]. Available: <http://arxiv.org/abs/1604.03635>
- [8] N. Nikhil and B. T. Morris, "Convolutional neural network for trajectory prediction," *Computer Vision - ECCV 2018 Workshops*, pp. 186–196, 2019. [Online]. Available: http://dx.doi.org/10.1007/978-3-030-11015-4_16
- [9] S. Jawed, E. Boumaiza, J. Grabocka, and L. Schmidt-Thieme, "Data-driven vehicle trajectory forecasting," 2019.
- [10] Y. Hu, W. Zhan, and M. Tomizuka, "Multi-modal probabilistic prediction of interactive behavior via an interpretable model," *CoRR*, vol. abs/1903.09381, 2019. [Online]. Available: <http://arxiv.org/abs/1903.09381>
- [11] B. Kim, C. M. Kang, J. Kim, S. H. Lee, C. C. Chung, and J. W. Choi, "Probabilistic vehicle trajectory prediction over occupancy grid map via recurrent neural network," *2017 IEEE 20th International Conference on Intelligent Transportation Systems (ITSC)*, Oct 2017. [Online]. Available: <http://dx.doi.org/10.1109/itsc.2017.8317943>
- [12] N. Deo and M. M. Trivedi, "Convolutional social pooling for vehicle trajectory prediction," *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, Jun 2018. [Online]. Available: <http://dx.doi.org/10.1109/cvprw.2018.00196>
- [13] S. H. Park, B. Kim, C. M. Kang, C. C. Chung, and J. W. Choi, "Sequence-to-sequence prediction of vehicle trajectory via lstm encoder-decoder architecture," *2018 IEEE Intelligent Vehicles Symposium (IV)*, Jun 2018. [Online]. Available: <http://dx.doi.org/10.1109/ivs.2018.8500658>
- [14] J. F. P. Kooij, N. Schneider, F. Flohr, and D. M. Gavrila, "Context-based pedestrian path prediction," in *Computer Vision – ECCV 2014*, D. Fleet, T. Pajdla, B. Schiele, and T. Tuytelaars, Eds. Cham: Springer International Publishing, 2014, pp. 618–633.
- [15] A. Alahi, K. Goel, V. Ramanathan, A. Robicquet, F. F. Li, and S. Savarese, "Social lstm: Human trajectory prediction in crowded spaces," 06 2016, pp. 961–971.
- [16] N. Lee, W. Choi, P. Vernaza, C. B. Choy, P. H. S. Torr, and M. Chandraker, "Desire: Distant future prediction in dynamic scenes with interacting agents," *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jul 2017. [Online]. Available: <http://dx.doi.org/10.1109/CVPR.2017.233>
- [17] G. Habibi, N. Jaipuria, and J. P. How, "Context-aware pedestrian motion prediction in urban intersections," 2018.
- [18] D. A. Ridet, N. Deo, D. Wolf, and M. M. Trivedi, "Understanding pedestrian-vehicle interactions with vehicle mounted vision: An lstm model and empirical analysis," 2019.
- [19] Y. Ma, X. Zhu, S. Zhang, R. Yang, W. Wang, and D. Manocha, "Trafficpredict: Trajectory prediction for heterogeneous traffic-agents," *ArXiv*, vol. abs/1811.02146, 2018.
- [20] A. H. Lang, S. Vora, H. Caesar, L. Zhou, J. Yang, and O. Beijbom, "Pointpillars: Fast encoders for object detection from point clouds," 2018.
- [21] T. Fernando, S. Denman, S. Sridharan, and C. Fookes, "Soft + hardwired attention: An lstm framework for human trajectory prediction and abnormal event detection," *Neural Networks*, vol. 108, p. 466478, Dec 2018. [Online]. Available: <http://dx.doi.org/10.1016/j.neunet.2018.09.002>
- [22] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," 2014.
- [23] A. Sadeghian, A. Alahi, and S. Savarese, "Tracking the untrackable: Learning to track multiple cues with long-term dependencies," *2017 IEEE International Conference on Computer Vision (ICCV)*, Oct 2017. [Online]. Available: <http://dx.doi.org/10.1109/ICCV.2017.41>
- [24] A. Gupta, J. Johnson, L. Fei-Fei, S. Savarese, and A. Alahi, "Social gan: Socially acceptable trajectories with generative adversarial networks," *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, Jun 2018. [Online]. Available: <http://dx.doi.org/10.1109/cvpr.2018.00240>
- [25] T. Zhao, Y. Xu, M. Monfort, W. Choi, C. Baker, Y. Zhao, Y. Wang, and Y. N. Wu, "Multi-agent tensor fusion for contextual trajectory prediction," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [26] T. Fernando, S. Denman, S. Sridharan, and C. Fookes, "Tracking by prediction: A deep generative model for multi-person localisation and tracking," *CoRR*, vol. abs/1803.03347, 2018. [Online]. Available: <http://arxiv.org/abs/1803.03347>