# Deep Learning for Self-Driving Vehicles

Safa Jameel Dawood Al-Kamil[1], Mohammed Salah Al-Radhi[2]
{safa.alkamil@stu.edu.iq[1], malradhi@tmit.bme.hu[2]}


Southern Technical University, Basra, Iraq[1]
Budapest University of Technology and Economics, Budapest, Hungary[2]

**Abstract.** Self-driving vehicles (SDV) and advanced safety features offering the greatest challenges and opportunities for Artificial Intelligence. The understanding of human intention is a very difficult task. As a result, predicting other drivers' future behaviour is critical for perceiving their past motion, analysing their interactions with other agents, and processing the data available from the scene. Automated driving systems (ADSs) promise to make driving safer, more comfortable, and more efficient. The Deep Structured Self-Driving Network (DSDNet) is proposed in this work that uses a single neural network to conduct object identification, motion prediction, and motion planning. The deep structured energy-based model, on the other hand, is improved. DSDNet also takes advantage of the expected forthcoming predicted actors to prepare a safe manoeuvre. Experiments The results reveal that it considerably increases detection, prediction, and planning performance.


**Keywords:** Self-driving vehicles, DSDNet, Deep learning, Prediction, and planning.

## 1 Introduction

Many discoveries in computer vision [1], robotics [2], and natural language processing (NLP) [3] were made possible by deep learning (DL) and artificial intelligence (AI). These methods have excellent control features on the future current autonomous driving in academics and industry. A self-driving vehicle (SDV) must monitor and predict surrounding actors' future behaviours, as well as plan safe manoeuvres. To be able to drive safely on the road. The success of DL, the prediction problem remains difficult due to the difficulties. Furthermore, motion planners must be improved in order to account for the uncertainty of forecasts. Parametric distributions have been used in the past to model multimodality in motion prediction. Because of their close-form inference, a combination of Gaussians [1, 2] is a natural approach. Nonetheless, deciding on the number of modes ahead of time is challenging. Furthermore, through training, these techniques suffer from mode collapse [3, 2, 4]. An alternative is to use neural networks to predict the data distribution through a series of training data. The multi – sensor operations carried out by CVAE [8], showed that the latent variables can be used to represent actor interactions, as shown in [5, 6, 7].

In this study, DSD- Net is suggested as a single neural network that leverages raw sensor data to jointly detect actors in the scene, forecast a multimodal distribution across their future behaviours, and build safe SDV plans. Typically, planning is framed as a problem of cost

minimization over trajectories. The fitness function can be manually constructed to ensure specific properties [16, 17,18, 19], or it can be trained from data via imitation learning or inverse reinforcement learning [20, 21]. These planners, however, believe that detection and prediction are exact and certain, which is not the case in fact. As a result, take into account the uncertainty of other players' actions and define collision avoidance in a probabilistic manner. Uncertainty-aware motion planning is explicitly applied to model the interactions between the SDV and the other dynamic agents, to achieve safer planning. Furthermore, our planning cost functions explicitly account for safety.

With a sample-based paradigm, this paper focuses on the costly probabilistic inference. Deep structured models, which use DNNs to provide the energy terms of probabilistic graphical models, have recently exploded in popularity as a way to encode prior knowledge (PGMs). DS models have been effectively used to several computer vision problems, such as semantic segmentation anomaly detection, and contour segmentation [22], by combining the powerful learning capability of DNNs with the task-specific structure imposed by PGMs. Inference for continuous random variables, on the other hand, is extremely difficult. We build a deep structured model called sample-based belief propagation (BP) that can learn complicated human behaviours from vast data processing based on past knowledge, which is influenced by this research. We also use a physically valid sampling approach to get over the issue of continuous variable inference.

## 2   Methods and experiments

In this study, we propose a new design technique (DSDNet) that discovers and predicts a generally regular multimodal allocation after temporal activities, as well as providing safe self-driving car movement plans.

The middle feature maps are evaluated using a backbone network before being developed further. Subsequently recognising performers using a detection header and taking into account their reactions, a deep organised probabilistic deduction module generates future trajectories, which indicates that allocations of actors will occur. Finally, the projected track takes into account the information stored in the feature maps, In addition, the model's likely futures. Our proposed strategy is depicted in Fig 1 as a high-level overview.
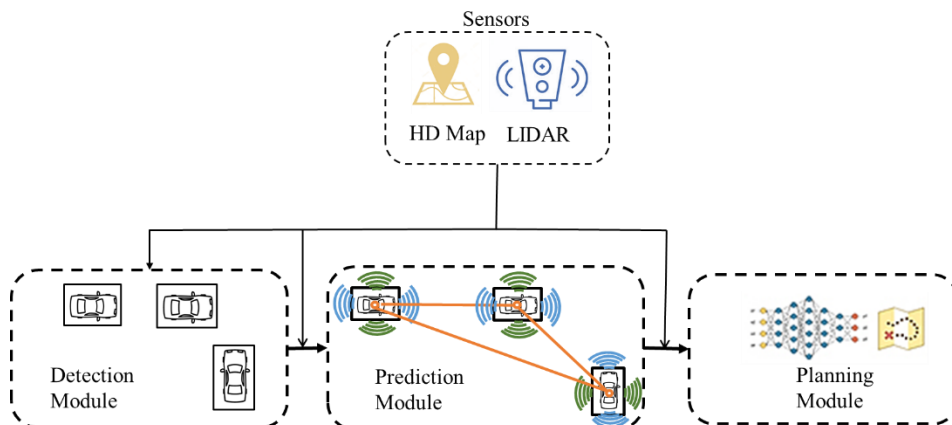
**Fig 1.** The multimodal trajectory prediction module is described in detail.

The following is a breakdown of the paper's structure: SDV's Methods and Experiments are explained in Section 2. The Experimental Evaluation is discussed in Section 3. And finally conclusions was addressed in section 4.

### 2.1 Object Detection and the Backbone Network

We vowelized the most recent 8 LiDAR sweeps to create a 3D tensor. LiDAR is a vital signal that can be sparse and utilised to detect and forecast the movement of performers. We utilise HD maps because they give us a good idea of what's going on. Besides the lanes are predicted by the traffic signal and like turning or straight into multiple pathways and let 3D LiDAR tensor to manage the control presentation. After that, we use a deep convolutional network backbone to process this 3D tensor and construct a backbone feature map, which must be followed by using a specific header pointer to the map the boundary squares in the scene. 2D convolution layers is applied, in which, one for organising location is taken by an performer and the other for retracting each actor's place offset, dimensions, direction, and velocity. The prediction and planning modules in this paper will choose the suitable prediction map for entry to provide both demeanors and safe planning.

### 2.2 Probabilistic Multimodal Trajectory Prediction

We used a path plan represented by a series of 2D intermediate sampled at T discrete timestamps on a birds-eye view (BEV). We calculate the movement prediction allocation and a movement of each sensor occurs, and time sample T.

### 2.3 Output Parameterization

We propose a solution from a limited number of samples to develop the required continuous space distribution. A few K samples to predict the future trajectory from random distribution for each actor, is illustrated in Fig. 2. Then limit each actor's conceivable future state to one of those K samples. The Neural Motion Planner (NMP) [9] is used. To maintain the whole diverse and reasonable distribution, a mix of circle, square, and colthood arcs are used.

### 2.4 Inference of Message Passing

Our motion planner demands us to examine all conceivable actor futures for safety reasons. As a result, motion forecasting that deduces the likelihood of each performer following a specific future path. As a result, we do marginal inference on the joint distribution. We also employ the crossing between the sum and product of data to assess the peripheral distribution of each fitness element, that considers the marginalisation impacts of the other elements under considerations.

### 2.5 Safe Motion Planning

Our last goal is to arrive our target and side-stepping collisions and following traffic laws in the motion planning module. To do this, we create fitness formula that reduce the complicity and to achieve excellent route compared to less safe other routes. The development idea was completed by determining the most cost-effective trajectory. We use a multi-task loss to train

the entire model (backbone, detection, prediction, and planning) in collaboration with a multi-class cost, which maintain the best supervision and the reliable training. We use a common detection methodology, that is a combination of classification and regression loss. We use cross-entropy between our discrete distribution and the genuine goal based on samples collected for the performers conduct. The main goal to achieve minimum Euclidian distances for the trajectory samples for the best predicable future trajectory.

This work established wide safety margins for potentially harmful behaviours such as collision paths. As a result, our approach can punish risky behaviours more severely.
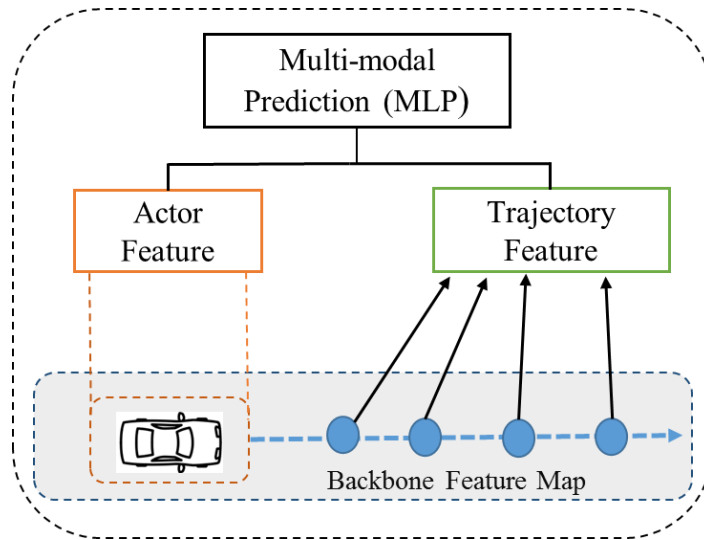


**Fig 2.** The neural header for evaluating Trajectory.

## 3  Experimental Evaluation

In this article, the model is tested for all three tasks: detection, prediction, and planning. On two large-scale real-world self-driving datasets, it outperforms the state-of-the-art on public benchmarks: scenes, as well as our own data ATG4D. In addition to the CARLA simulation dataset prediction module is tested. The benefits of explicitly modelling actor interactions are demonstrated in this work. When compared to rival systems, this article planning module produces the safest planning outcomes and reduces the rate of collisions and lane violations significantly. The proposed model continues the detection even over a single backbone with implication. In the supplementary material, we detail the dataset as well as the implementation.
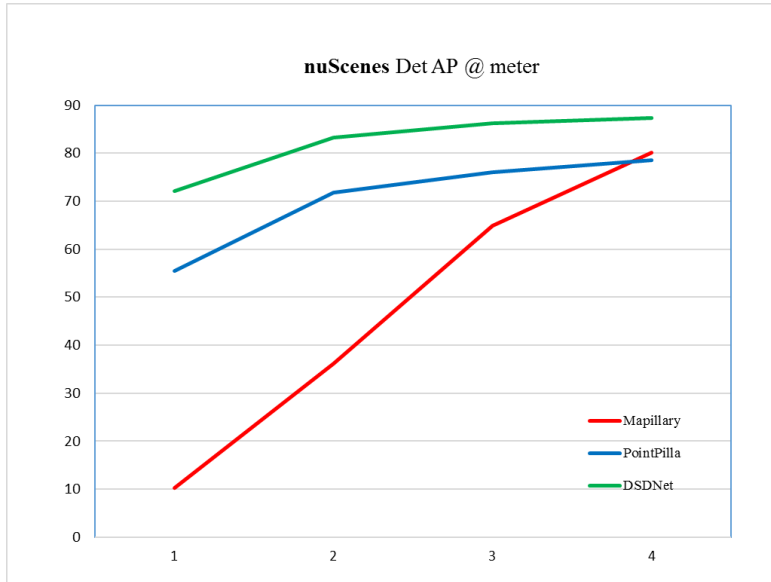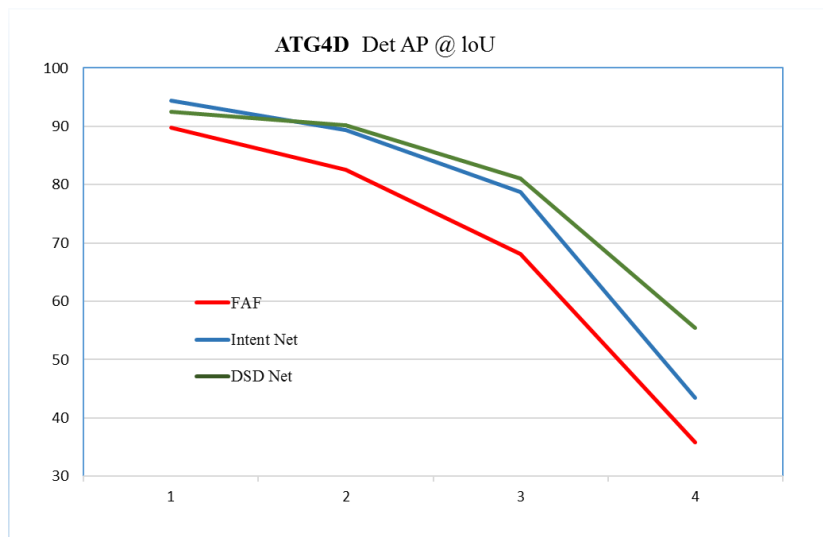
**Fig 3.** Detection performance of scenes.



**Fig 4.** The results of ATG4D's detection performance.

On both datasets, as demonstrated in Figs 3 and 4, our technique produces the best results. This is significant because conventional baselines employ L2 as a training target and are thus directly favoured by the L2 error metric, whereas our techniques learn correct distributions and capture multi-modality via cross-entropy loss. Multimodal approaches are regarded to have the lowest score in this metric. We show that modelling multi-modality while attaining lower L2 error is

attainable due to the adaptable performers' behaviour. The method lowers collisions between the anticipated trajectories of the actors, demonstrating the value of our multi-agent interaction modelling. This evaluation is carried out on CARLA, which has been used in all prior techniques. As seen in Fig 4, our strategy exceeds prior best findings by a wide margin. Our detections, forecasts, the expected uncertainties that are all visualised. We use different colours for various future timestamps to denote high-probability performers' future positions estimated via our prediction module. The predictions are certain when vehicles drive along the lanes (left), but when vehicles approach a junction, multi-modal predictions are seen (middle, right).

## 4  Conclusion

Self-driving vehicles (SDV) and how the safety feature offering the highest challenge for Artificial Intelligence has been presented. Our planning modules produce the safest planning results, with significant reductions in the rate of collisions and lane violations. In this paper, we have suggested DSDNet, to serve under a unified framework. Perception, prediction, and planning.

Experiments Results on a self-driving dataset show that our model improves the performance of detection, prediction, and planning remarkably. We want to understand more about structured energies and how to build excellent reward functions for autonomous agents, as well as how to incorporate safety into decision-making reinforcement learning systems.

## References

[1]  Chai Y., Sapp B., Bansal M., Anguelov D. Multipath: Multiple probabilistic anchor trajectory hypotheses for behavior prediction. arXiv preprint arXiv:1910.05449 (2019).

[2]  Hong J., Sapp B., Philbin J. Rules of the road: Predicting driving behavior with a convolutional model of semantic interactions. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 8454-8462 (2019).

[3]  Jain A., Casas S., Liao R., Xiong Y., Feng S., Segal S., Urtasun R. Discrete residual flow for probabilistic pedestrian behavior prediction. arXiv preprint arXiv:1910.08041 (2019).

[4]  Rhinehart N., Kitani K.M., Vernaza P. R2p2: A reparameterized pushforward policy for diverse, precise generative path forecasting. In: European Conference on Computer Vision. pp. 794-811. Springer, Cham (2018).

[5]  Rhinehart N., McAllister R., Kitani K., Levine S. Precog: Prediction conditioned on goals in visual multi-agent settings. arXiv preprint arXiv:1905.01296 (2019).

[6]  Lee N., Choi W., Vernaza P., Choy C.B., Torr P.H., Chandraker M. Desire: Distant future prediction in dynamic scenes with interacting agents. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 336{345 (2017).

[7]  Tang Y.C., Salakhutdinov R. Multiple futures prediction. arXiv preprint arXiv:1911.00997 (2019).

[8]  Sohn K., Lee H., Yan X. Learning structured output representation using deep conditional generative models. In: Advances in neural information processing systems. pp. 3483-3491 (2015).

[9]  Zeng W., Luo W., Suo S., Sadat A., Yang B., Casas S., Urtasun R. End-toend interpretable neural motion planner. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 8660{8669 (2019).

[10] Liang M., Yang B., Zeng W., Chen Y., Hu R., Casas S., Urtasun R. Pnpnet: End-to-end perception and prediction with tracking in the loop. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 11553-11562 (2020).

[11] Wang T.H., Manivasagam S., Liang M., Yang B., Zeng W., Raquel U. V2vnet: Vehicle-to-vehicle communication for joint perception and prediction. In: Proceedings of the European Conference on Computer Vision (ECCV) (2020).

[12] Luo W., Yang B., Urtasun R. Fast and furious: Real time end-to-end 3d detection, tracking and motion forecasting with a single convolutional net (2015).

[13] Casas S., Luo W., Urtasun R. Intentnet: Learning to predict intention from raw sensor data. In: Billard, A., Dragan, A., Peters, J., Morimoto, J. (eds.) Proceedings of The 2nd Conference on Robot Learning. Proceedings of Machine Learning Research, vol. 87, pp. 947-956. PMLR (29{31 Oct 2018).

[14] Casas S., Gulino C., Suo S., Luo K., Liao R., Urtasun R. Implicit latent variable model for scene-consistent motion forecasting. In: Proceedings of the European Conference on Computer Vision (ECCV) (2020).

[15] Casas S., Gulino C., Suo S., Urtasun R. The importance of prior knowledge in precise multimodal prediction. In: IROS (2020).

[16] Buehler M., Iagnemma K., Singh S. The DARPA urban challenge: autonomous vehicles in city traffi c, vol. 56. springer (2009).

[17] Fan H., Zhu F., Liu C., Zhang L., Zhuang L., Li D., Zhu W., Hu J., Li H., Kong Q. Baidu apollo em motion planner. arXiv preprint arXiv:1807.08048(2018).

[18] Montemerlo M., Becker J., Bhat S., Dahlkamp H., Dolgov D., Ettinger S., Haehnel D., Hilden T., Homann G., Huhnke B., et al. Junior: The Stanford entry in the urban challenge. Journal of eld Robotics 25(9), 569{597 (2008).

[19] Ziegler J., Bender P., Dang T., Stiller C. Trajectory planning for bertha|a local, continuous method. In: Intelligent Vehicles Symposium Proceedings, 2014 IEEE. pp. 450-457. IEEE (2014).

[20] Sadat A., Ren M., Pokrovsky A., Lin Y.C., Yumer E., Urtasun R. Jointly learnable behavior and trajectory planning for self-driving vehicles. arXiv preprint arXiv:1910.04586 (2019).

[21] Wulfmeier M., Ondruska P., Posner I. Maximum entropy deep inverse reinforcement learning. arXiv preprint arXiv:1507.04888 (2015).

[22] Marcos D., Tuia D., Kellenberger B., Zhang L., Bai M., Liao R., Urtasun R. Learning deep structured active contours end-to-end. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (2018).