

A Recurrent Neural Network Solution for Predicting Driver Intention at Unsignalized Intersections

Alex Zyner, Stewart Worrall, and Eduardo Nebot¹

Abstract—In this paper we present a system capable of inferring intent from observed vehicles traversing an unsignalized intersection, a task critical for the safe driving of autonomous vehicles, and beneficial for advanced driver assistance systems (ADAS). We present a prediction method based on Recurrent Neural Networks (RNNs) that takes data from a Lidar based tracking system similar to those expected in future smart vehicles. The model is validated on a roundabout, a popular style of unsignalized intersection in urban areas. We also present a very large naturalistic dataset recorded in a typical intersection during two days of operation. This comprehensive dataset is used to demonstrate the performance of the algorithm introduced in this work. The system produces excellent results, giving a significant, 1.3 second prediction window before any potential conflict occurs.

Index Terms—Intelligent Transportation Systems, Deep Learning in Robotics and Automation

I. INTRODUCTION

A. Motivation

NAVIGATING intersections that do not have traffic signals is a complex task that requires significant interaction and prediction of other drivers in the scene. While an experienced driver may give such a task little thought, it is still a complex problem for an intelligent vehicle, and as such is an open area of research. One such intersection style is the roundabout, which is widespread in many urban areas as they do not require expensive traffic lights, and have higher throughput than a four-way stop. In this paper we present an algorithm that can predict the movement of an observed vehicle, given a short segment of tracking data from a intelligent vehicle's onboard sensors. This ability to predict the intention of drivers at an intersection is critical for the safe driving of autonomous vehicles, and beneficial for advanced driver assistance systems.

B. Problem Definition

The roundabout is commonly used in Australia as an unsignalized intersection due to its high throughput in medium to low traffic areas. Australian roundabouts have a number

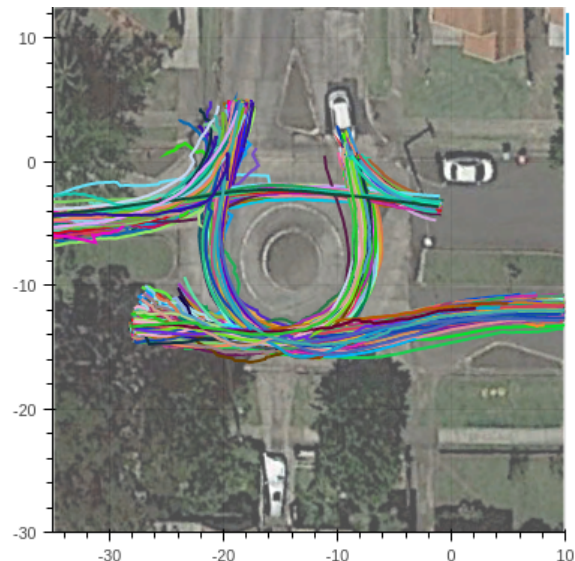


Fig. 1. An overlay of 1000 of the recorded tracks at the intersection studied.

of design features that make them very different from European roundabouts. Australian roundabouts are tangential, as opposed to European which are radial. These differences can be seen in Figure 2. Tangential roundabouts encourage speed, and are coupled with good visibility. This encourages faster travel, which equates to higher throughput. An artifact of this is that the negotiation between drivers is mostly done during the approach. As such, many drivers will not slow their speed to assert their ordering in the roundabout, even though it may be contrary to the road rules. An experienced human driver will recognize these traits and be able to accommodate them, and so safely navigate the intersection. This work aims to emulate this human intuition, by presenting a system that is able to recognize the intention of observed vehicles in an intersection, validated on a very large, naturalistic dataset.

C. Contributions

Future vehicle ADAS and safety systems will incorporate contextual information, such as driver intention and expected behaviour to achieve anticipatory driving. Gathering this information is a challenging task because it cannot be measured directly, it must be inferred from sensor data. This paper presents an innovative algorithm to provide an estimation of driver intention at an unsignalized roundabout.

This inference is to be used as an additional source of information to improve the safety of the vehicle's operation.

Manuscript received: September, 10th, 2017; Revised November, 28th, 2017; Accepted January, 22nd, 2018.

This paper was recommended for publication by Wan Kyun Chung upon evaluation of the Associate Editor and Reviewers' comments. This work was supported by: The Australian Government through the Australian Research Council Discovery Grant DP160104081, and the Next Generation Vehicle project, funded by University of Michigan, Ford Motor Company

¹Authors are with the Australian Centre for Field Robotics (ACFR) at the University of Sydney (NSW, Australia). {a.zyner, s.worrall, e.nebot} (at) acfr.usyd.edu.au

Digital Object Identifier (DOI): see top of this page.

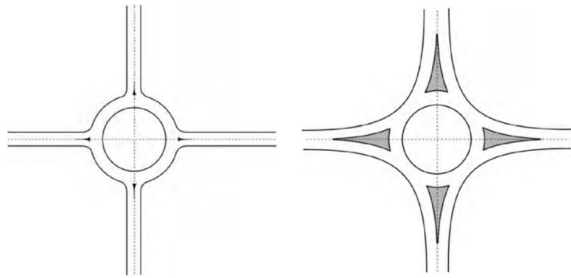


Fig. 2. A European style radial (left) and an Australian tangential (right) roundabout [3]. The tangential roundabout encourages speed as the driver's path is straighter than the radial roundabout. Tangential roundabouts are typically coupled with high visibility during the approaches.

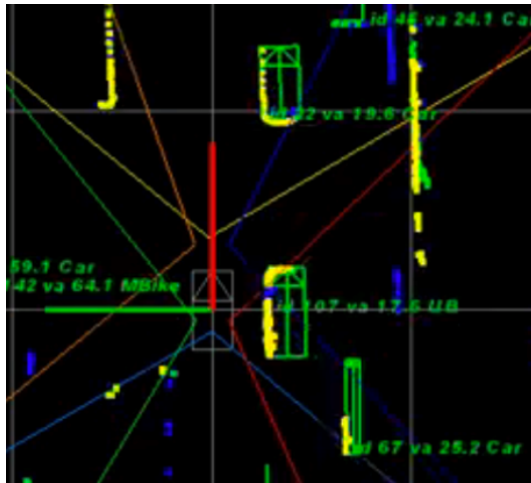


Fig. 3. Sample data collected from the Ibeo Capture System. Each of the 6 Lidars can be seen surrounding the vehicle, with their 110 degree field of view highlighted in a coloured 'V'. Point returns are labeled in green, yellow, blue and red depending on which line of the scan detected a hit. Tracked vehicles are labeled in green.

Ideally, we would like to make decisions based on multiple sources of information to achieve a high integrity solution. This concept has been referred to by the authors as system integrity, and has been successfully applied to vehicle control systems [1] as part of large scale deployments of autonomous field robotics applications [2]. The fundamental concept of integrity is based on using multiple sources of information obtained with sensors that rely on different physical principles. The contribution of the work presented in this paper is to provide one source of information for estimating driver intent. It is expected that this will be used in combination with other information sources in safety and autonomous driving applications.

The paper is organized as follows. Section II presents related work. Section III presents the dataset, describing how it was collected, and presents a summary of the data. A RNN architecture for intention prediction is presented in Section IV. The experimental procedure is presented in V and the results in Section VI. We present final remarks in the conclusion, Section VII.



Fig. 4. The data collection vehicle, parked at the roundabout studied.

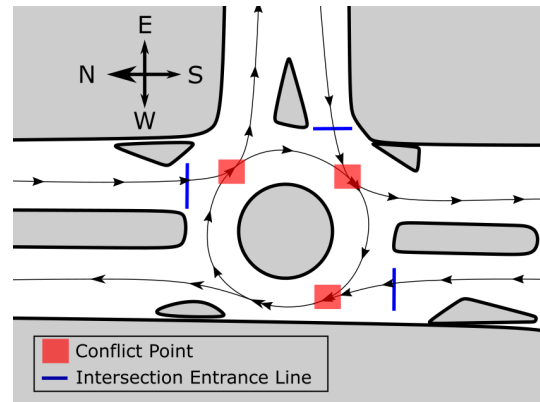


Fig. 5. Diagram of the intersection studied. Note, as this is a left hand drive road, so vehicles traverse the roundabout in a clockwise fashion. Conflict points, the points of collision between a vehicle on the roundabout, and a vehicle entering the roundabout, are marked with a red square.

II. RELATED WORK

Predicting the intention of road users around an intersection is a widely studied problem. Techniques such as Hidden Markov Models [4], [5], RNNs [6] or Support Vector Machines [7] have been proven successful for use in intention prediction, however most of the existing work generally rely on tracking data taken from the ego vehicle. Data such as GPS, steering wheel encoding, [8], [9] or even internal cabin cameras to study driver's gaze and head position [10], [11] is recorded to use as leading indicators to predict the next maneuver of the vehicle.

These works rely on vehicles being outfitted with such a detection system, so that the vehicles may then broadcast the driver's intent via vehicle to vehicle communication, or otherwise act on the sole information of the ego vehicle only. As there will always be vehicles without this type of technology on the road for the foreseeable future, it is necessary for an intelligent vehicle to infer intent about cars external to the ego vehicle. This data could be collected by a smart intersection, using overhead cameras or Lidars [12]. A Long-Short Term Memory (LSTM) based solution is presented by Phillips et al. [13] that focuses on large, multi lane intersections in the US., using features such as speed, which lane the car is in and information about the surrounding six vehicles.

An intelligent vehicle will not be able to solely rely on smart infrastructure to gather information, the system must be able to infer intent of surrounding vehicles using on-board sensors

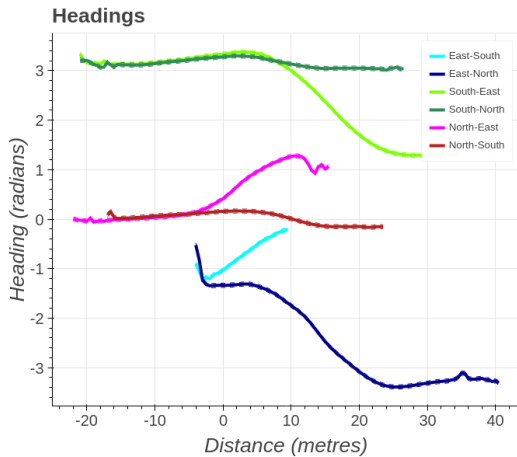


Fig. 6. Headings profile of all tracks, grouped into the six classes in 'origin-destination' pairs. The x axis is distance traveled relative to the intersection entrance. Mean heading is plotted as a solid line. Standard deviation is not visible in this graph due to its scale.

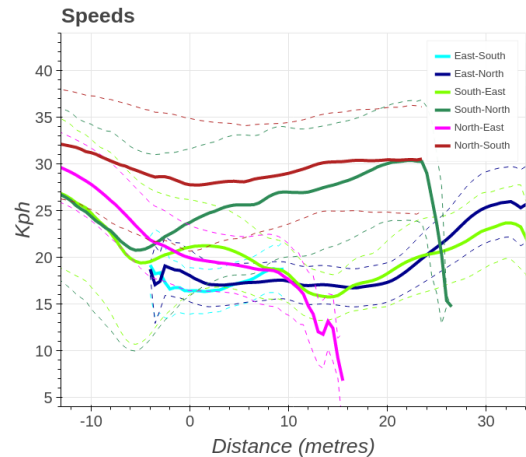


Fig. 7. Speed profile of all tracks, grouped into the six classes in 'origin-destination' pairs. The x axis is distance traveled relative to the intersection entrance. Mean speed is plotted as a solid line. One standard deviation is depicted as a dashed line.

alone. Muffert et al. [14] present a stixel based stereo vision solution at a large urban roundabout, using a time-to-collision based metric. A similar method is used by Barth et al. [15] focusing on a turn-across-lane scenario using both real and synthetic data. Using the KITTI dataset, Khosroshahi et al. [16] present an LSTM based method to classify tracks through a signalized intersection, as recorded by Lidar data. The data used is limited as only 49 tracks were used, and they do not present a predictive model. A further review of recent vehicle prediction algorithms is presented by Lefevre et al. [17].

III. DATASET

A. Data Collection Vehicle

The vehicle used for data collection is outfitted with a ibeo.HAD Feature Fusion detection and tracking system. This system uses 6 ibeo LUX 4 beam, 25 Hz Lidar scanners to identify road users at a range of up to 200m, and has an on-board computer for classification and tracking, in real time. A sample output of this system can be seen in Figure 3. Data about each detected vehicle is collected, including X/Y relative positioning (metres), velocity (metres/second), heading (radians), size (width/height metres), classification [bike, car, heavy vehicle, pedestrian], and classification confidence. Indicator status is not collected, or recorded. This data is recorded at a rate of 25 Hz. This system is similar to systems an autonomous vehicle is expected to have.

B. The Intersection

The data was collected by the vehicle parked in the location demonstrated in Figure 4. This allows for a clear view of approaching vehicles from all directions, and also simulates the perspective of a vehicle approaching from this direction. A diagram of this three-way roundabout can be seen in Figure 5. This roundabout exists in Australia, where vehicles drive on the left hand side of the road, and traverse a roundabout in a clockwise direction. The data was grouped by the track's origin and destination.

TABLE I
SUMMARY OF DATA COLLECTED AT THE ROUNDABOUT, GROUPED BY ORIGIN AND DESTINATION CLASSES.

Origin	Destination			
	East	North	South	Total
East	0	2588	385	2973
North	2705	0	607	3312
South	1230	777	0	2007
Total	3935	3365	992	8292

TABLE II
TABLE OF THE AVERAGE SPEEDS OF EACH CLASS, AND THEIR TRANSITION THROUGH THE INTERSECTION

Origin-Destination	Average Speed (kph)	Std. Dev (Kph)	Turn
East-South	17.20	2.54	Left
East-North	18.99	3.89	Right
North-South	29.54	6.15	Straight
North-East	22.32	5.40	Left
South-North	25.87	8.46	Straight
South-East	20.43	5.91	Right

These tracks were recorded over 2 days, and consist of 20 hours of natural traffic passing through the area. This has resulted in the collection of over 8000 vehicle tracks passing through the intersection. To the author's knowledge, this is the largest dataset collected from onboard vehicle sensors for an intersection study.

The overall class distribution can be seen in Table I. The collation of all tracks recorded is visible in Figure 1, where the car is parked at position (0,0). Due to how the car was parked on the side of the road, some of the vehicles being tracked were occluded by another vehicle leaving the intersection. This makes labeling difficult, as the vehicle needs to be observed leaving the intersection to properly label the training data. For this reason, partial tracks were discarded for training. At runtime, the system only needs a fraction of a recording to perform classification, and thus the network can predict the car well before it gets to its destination.

The labeling process was implemented as follows. The three

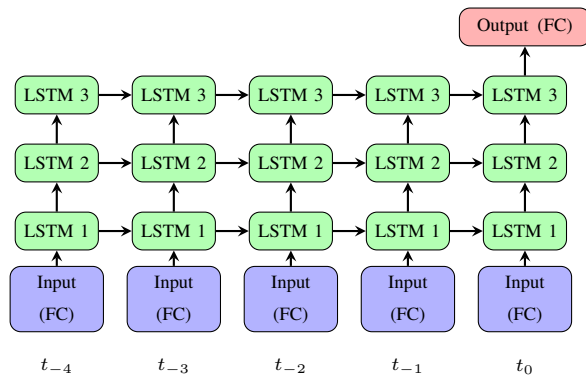


Fig. 8. Diagram representing a RNN of example length 5, that is given data for time t . Here the model is input with a sequence of length 5, that ends at the nominated time-step t , and outputs a class in the set [East, North, South]. The first layer consists of a single fully connected (FC) layer. The next layers are three recurrent layers represented in green. Here the horizontal links represent the hidden layer transition between timesteps of an LSTM.

exits, and three entrances to the intersection were hand labeled in the (stationary) vehicle reference frame. A vehicle track is allocated a category if the track passes through both an entrance, and an exit. The distance traveled at each point from (or to) the intersection entrance is calculated, as it is used as a result metric. The total count for each class is presented in Table I. Speed profiles for each class are presented in Table II.

C. Dataset Findings

The average speed and headings for each class can be seen in Figures 6 and 7. Comparing the heading of the two track classes that start in the north, they begin to visibly diverge within 2m of the intersection entrance (denoted as the 0m mark). The speeds diverge as well, however it takes significant distance for the one standard deviation line of each class to diverge, which occurs at the 6 metre mark. As these two classes are a short, left turn, and a straight ahead, it is expected that tracks in these two classes diverge very early into the intersection. Track classes that begin in the East also diverge quite quickly in heading, as this is comparing a short, left turn and a longer right turn. The speed profiles of each class never diverge past one standard deviation. Finally, the tracks that begin in the South are the most ambiguous case, as these classes take significant distance to become statistically deviant from one another. After about 11 metres of travel, the headings separate, and the speeds also become distinct from one another.

This dataset can be downloaded at the following website: <http://its.acfr.usyd.edu.au/datasets/>

IV. NETWORK ARCHITECTURE

In this paper, we present a model architecture used to infer prediction of this data, which is based on an RNN. These networks are especially useful when dealing with time series data, and have been shown to work in areas such as pedestrian path prediction [18], and text parsing [19]. The recurrent nature of the network allows the system to consider both past and present data. It does this by having a copy of the network

for each time-step, and passing internal activations to forward time-steps.

The RNN model presented is used to interpret time series data about an externally observed vehicle. We do this by having one recurrence of the network per time-step, and after a chosen number of time-steps, we allow the network to make a prediction. A diagram of this network format can be seen in Figure 8.

The input features of the model are X / Y position relative to the vehicle, heading in radians, and speed in metres per second. The data is normalized over the entire training set, and then input directly into the network. The first layer of the network is a fully connected layer, followed by three layers deep of recurrent layers. The recurrent cell chosen for this network is a LSTM with peephole connections [20]. A dropout value of 0.5 was used on the inter-layer connections only, and not the recurrent connections [21]. After applying a softmax classifier, the output of the network is the three destinations of the intersection: East, North and South. The overall goal of the model is to correctly predict the destination as early as possible.

In order to properly train and evaluate a network designed for parsing sequential fragments of time series data, the dataset is preprocessed in the following manner. Each track in the collection is split into every possible sequential sequence that can be input into the RNN. This split is described in the following two equations. In equation (1), M is the set of all maneuver samples, T is the length of each the j^{th} track in M . x_t is the recorded data for the j^{th} track at time t , and y is the destination of the j^{th} track. To train the model, all tracks are split up into every possible consecutive sequence of length k , where k is defined as the number of steps in the RNN. This split is defined in equation (2). Here S is the complete preprocessed data set, where w consecutive training samples are taken from each track j . For model evaluation and training, the train / validation / test split is done track-wise, to ensure no training data overlaps the test or validation set.

$$M = \{(x_1, x_2, x_3, \dots, x_T)_j, y_j\}_{j=1}^N \quad (1)$$

$$S = \{ \{(x_w, x_{w+1}, x_{w+2}, \dots, x_{w+k}), y_j\}_{w=1}^{T-k} \}_{j=1}^N \quad (2)$$

In order to properly score the network, the tracks are aligned by distance traveled from the entrance of the intersection. For each track, the relative distance traveled before or after crossing the line at the entrance to the roundabout is used. As each data sample contains multiple time steps, the furthest distance in the set is used. This results in a fair comparison between RNNs of differing lengths.

V. EXPERIMENTS

The network used in the analysis has the following parameters. Three recurrent layers are used, each of 512 nodes width. A single dense layer is used as an input layer, of width 256. ADADELTA [22] training was used, with a learning rate of 0.03. We split the dataset into 4560 tracks for training, 1658 for validation of hyperparameters, and 2074 for final testing. Three

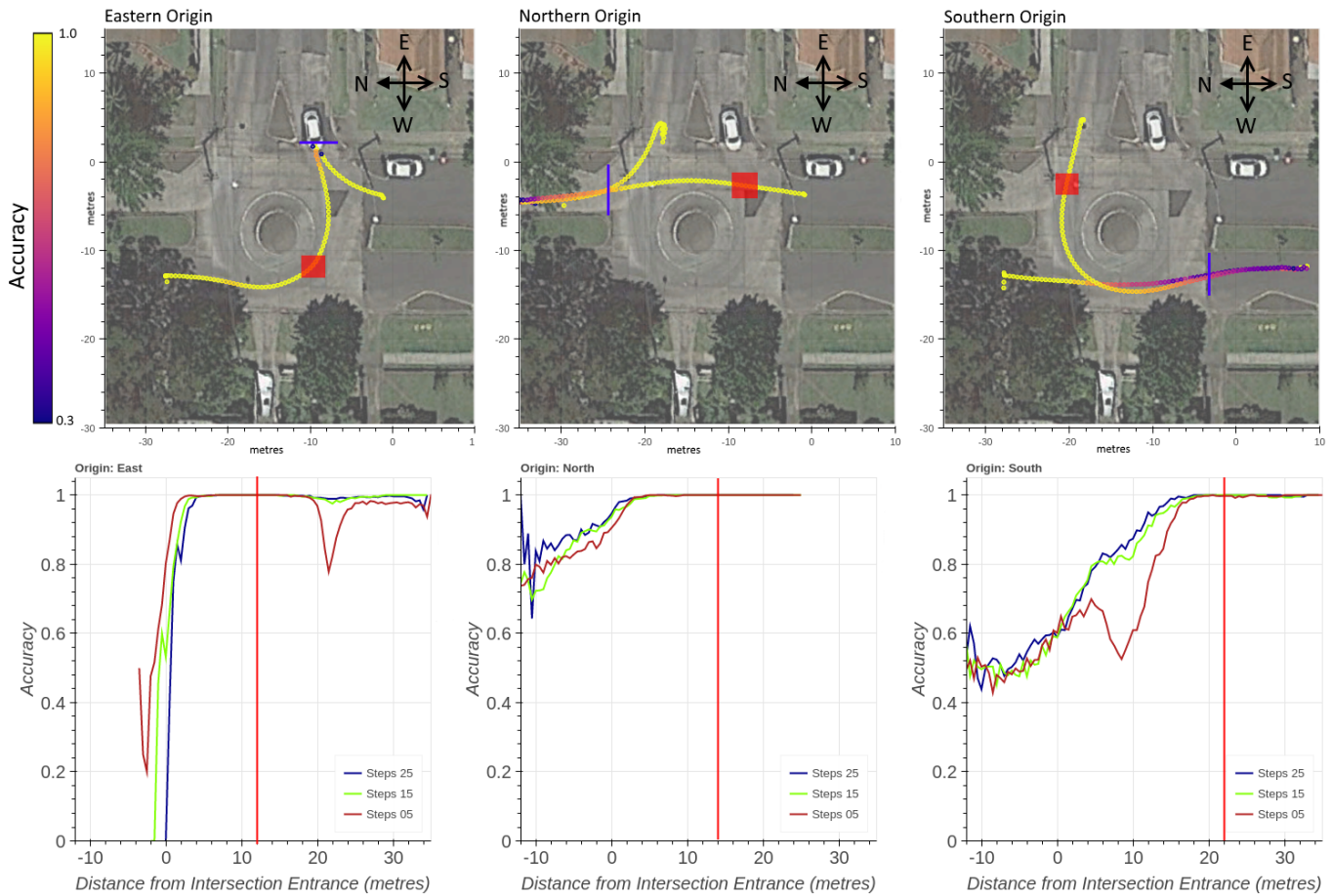


Fig. 9. Results collated by the approach road each driver took towards the intersection. The upper three figures are plots of the average track traversed for that origin-destination pair. The tracks are colored by accuracy of the 15 step RNN, the lighter (or more yellow) the color, the more accurate the result. A red square indicates the appropriate conflict point. The lower three figures are plots of accuracy vs distance traveled relative to the start of the intersection. Each of the three different lengths are plotted in these figures. The red line on these graphs indicate the position of the conflict point for that origin.

network lengths were chosen for investigation: 5, 15, and 25 time-steps long, which correlate to 0.2, 0.6 and 1 seconds of data, respectively. These networks are tested completely independent of each other. Each network was trained on a single Nvidia Geforce 1080 GPU, and no network took longer than 6 hours to train. After training, running inference on the network takes 60ms, which makes it feasible for real-time deployment. The network was written in Tensorflow [23].

VI. RESULTS

The graphs in Figure 9 display the results of the algorithm used to predict the destination of a vehicle on the roundabout. The three lower line charts are split by origin of the vehicle, and the x axis represented distance traveled to or from the intersection entrance. Here, results from RNNs of lengths 5, 15 and 25 are displayed. Plotting the accuracy relative to the distance traveled was chosen as it readily demonstrates the gradient of difficulty of the problem. A car that is far from approaching the intersection gives little to no indication of the destination, while it is trivial to determine the destination of a car that is already at its destination.

The three upper topological charts display these results for the RNN with a length of 15 steps. Here the data is displayed as an overlay on the map, to better visualize where the conflict points, and accuracies lie. The results are displayed as a color on the average path for that origin/destination pair, with yellow corresponding to high accuracy, and blue low accuracy. The conflict point for that particular entrance, which is the point of potential collision between traffic on the roundabout, and traffic incoming to the roundabout, is displayed in red. Overall the network behaves exceedingly well, having an excellent classification accuracy well before any conflict points in the intersection.

A. Eastern Origin

Here the vehicles are either making a close, left turn or a long right turn. The heading profile of each vehicle path makes it immediately obvious as to the vehicle's destination, and all three classifiers reach 99% at 4 metres traveled, giving about 1.6 seconds lead time before reaching a conflict point when considering the average speed travelled. The RNN of length 5 achieves this accuracy slightly earlier. This is because the

shorter RNN needs less data to produce a result, and so by the time the vehicle was tracked by the ego vehicle's sensors, it was fairly clear where the vehicle was going. What makes this set particularly interesting is the fact that the shorter RNN classifier loses accuracy at around 22 metres distance. This is because the vehicle's profile closely matches that of a south-to-east traveler, especially in speed.

Longer observations do not suffer from this shortfall, as the network is able to remember past positions, and so it would have better information about the vehicle's origin.

B. Northern Origin

Cars traveling from the north are either making a short, left turn, or are continuing straight. Drivers passing straight through the intersection can maintain speed through the intersection. The networks can easily pick up on this, with all three giving 95% accuracy at the intersection entrance and 99% accuracy at 5 metres distance, well before the nearest conflict point at a distance of 14 metres. Given the average speed, this is approximately a 1.3 second lead time before any potential collision.

C. Southern Origin

Vehicles approaching from the south may either continue straight, or make a large right turn. The longer classifiers here show much greater accuracy earlier on as compared to the RNN of length 5. The 15 and 25 time step classifiers converge to an accuracy of 99% at 16 metres traveled, which is when the vehicle would pass the first exit. This is well before the nearest contact point at approximately 22 metres distance, giving around a 1.3 second lead time when considering the average speed of the vehicles.

D. Overall Performance

The system as a whole produces excellent results, with all classifiers producing a very high accuracy well before any potential conflict point. There is evidence that networks fed with more history perform better, but there is diminished returns after about 0.6 seconds, which correlates to a network length of 15 time-steps for this system. That is to say, the system produces its best results with only a 0.6 second track of a vehicle, allowing for a vehicle's intention to be recognized with only a very short observation.

VII. CONCLUSION

In this paper we have presented a method for predicting the behaviour of naturalistic vehicles at a urban, single lane roundabout. We validate this algorithm on a very large and unique dataset, that is collected with a vehicle perception system that is similar to those expected on future smart vehicles. We demonstrate that this model achieves practical results, giving a significant, 1.3 second prediction window before any potential conflict, and only needs a short (0.6 second) amount of observation data. This shows that this system may be feasibly deployed in an autonomous vehicle or appropriately equipped ADAS vehicles.

REFERENCES

- [1] E. M. Nebot and H. Durrant-Whyte, "A high integrity navigation architecture for outdoor autonomous vehicles," *Robotics and Autonomous Systems*, vol. 26, no. 2-3, pp. 81–97, 1999.
- [2] S. Sukkarieh, E. M. Nebot, and H. F. Durrant-Whyte, "A high integrity imu/gps navigation loop for autonomous land vehicle applications," *IEEE Transactions on Robotics and Automation*, vol. 15, no. 3, pp. 572–578, 1999.
- [3] A. Wilke, J. Lieswyn, and C. Munro, "Assessment of the effectiveness of on-road bicycle lanes at roundabouts in australia and new zealand," Tech. Rep., 2014.
- [4] H. Berndt and K. Dietmayer, "Driver intention inference with vehicle onboard sensors," in *Vehicular Electronics and Safety (ICVES), 2009 IEEE International Conference on*. IEEE, 2009, pp. 102–107.
- [5] T. Streubel and K. H. Hoffmann, "Prediction of driver intended path at intersections," in *2014 IEEE Intelligent Vehicles Symposium Proceedings*. IEEE, 2014, pp. 134–139.
- [6] A. Zyner, S. Worrall, J. Ward, and E. Nebot, "Long short term memory for driver intent prediction," in *Intelligent Vehicles Symposium (IV), 2017 IEEE*. IEEE, 2017, pp. 1484–1489.
- [7] P. Kumar, M. Perrollaz, S. Lefevre, and C. Laugier, "Learning-based approach for online lane change intention prediction," in *Intelligent Vehicles Symposium (IV), 2013 IEEE*. IEEE, 2013, pp. 797–802.
- [8] A. Mudgal, S. Hallmark, A. Carriquiry, and K. Gkritza, "Driving behavior at a roundabout: A hierarchical bayesian regression analysis," *Transportation research part D: transport and environment*, 2014.
- [9] M. Zhao, D. Kathner, M. Jipp, D. Soffker, and K. Lemmer, "Modeling driver behavior at roundabouts: Results from a field study," in *Intelligent Vehicles Symposium (IV), 2017 IEEE*. IEEE, 2017.
- [10] E. Ohn-Bar, A. Tawari, S. Martin, and M. M. Trivedi, "On surveillance for safety critical events: In-vehicle video networks for predictive driver assistance systems," *Computer Vision and Image Understanding*, vol. 134, pp. 130–140, 2015.
- [11] A. Jain, A. Singh, H. S. Koppula, S. Soh, and A. Saxena, "Recurrent neural networks for driver activity anticipation via sensory-fusion architecture," in *2016 IEEE International Conference on Robotics and Automation (ICRA)*, May 2016, pp. 3118–3125.
- [12] E. Strigel, D. Meissner, F. Seeliger, B. Wilking, and K. Dietmayer, "The ko-per intersection laserscanner and video dataset," in *Intelligent Transportation Systems (ITSC), 2014 IEEE 17th International Conference on*. IEEE, 2014, pp. 1900–1901.
- [13] D. J. Phillips, T. A. Wheeler, and M. J. Kochenderfer, "Generalizable intention prediction of human drivers at intersections," in *Intelligent Vehicles Symposium (IV), 2017 IEEE*. IEEE, 2017, pp. 1665–1670.
- [14] M. Muffert, T. Milbich, D. Pfeiffer, and U. Franke, "May i enter the roundabout? a time-to-contact computation based on stereo-vision," in *Intelligent Vehicles Symposium (IV), 2012 IEEE*. IEEE, 2012.
- [15] A. Barth and U. Franke, "Tracking oncoming and turning vehicles at intersections," in *Intelligent Transportation Systems (ITSC), 2010 13th International IEEE Conference on*. IEEE, 2010, pp. 861–868.
- [16] A. Khosroshahi, E. Ohn-Bar, and M. M. Trivedi, "Surround vehicles trajectory analysis with recurrent neural networks," in *Intelligent Transportation Systems (ITSC), 2016 IEEE 19th International Conference on*. IEEE, 2016, pp. 2267–2272.
- [17] S. Lefevre, 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.
- [18] 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, pp. 961–971.
- [19] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," *CoRR*, vol. abs/1301.3781, 2013. [Online]. Available: <http://arxiv.org/abs/1301.3781>
- [20] H. Sak, A. W. Senior, and F. Beaufays, "Long short-term memory recurrent neural network architectures for large scale acoustic modeling," in *INTERSPEECH*, 2014, pp. 338–342.
- [21] W. Zaremba, I. Sutskever, and O. Vinyals, "Recurrent neural network regularization," *arXiv preprint arXiv:1409.2329*, 2014.
- [22] M. D. Zeiler, "Adadelata: an adaptive learning rate method," *arXiv preprint arXiv:1212.5701*, 2012.
- [23] "TensorFlow: Large-scale machine learning on heterogeneous systems," 2015, software available from tensorflow.org. [Online]. Available: <http://tensorflow.org/>