

MANTRA: Memory Augmented Networks for Multiple Trajectory Prediction

Francesco Marchetti Federico Becattini Lorenzo Seidenari Alberto Del Bimbo
MICC, University of Florence

name.surname@unifi.it

Abstract

Autonomous vehicles are expected to drive in complex scenarios with several independent non cooperating agents. Path planning for safely navigating in such environments can not just rely on perceiving present location and motion of other agents. It requires instead to predict such variables in a far enough future. In this paper we address the problem of multimodal trajectory prediction exploiting a Memory Augmented Neural Network. Our method learns past and future trajectory embeddings using recurrent neural networks and exploits an associative external memory to store and retrieve such embeddings. Trajectory prediction is then performed by decoding in-memory future encodings conditioned with the observed past. We incorporate scene knowledge in the decoding state by learning a CNN on top of semantic scene maps. Memory growth is limited by learning a writing controller based on the predictive capability of existing embeddings. We show that our method is able to natively perform multi-modal trajectory prediction obtaining state-of-the-art results on three datasets. Moreover, thanks to the non-parametric nature of the memory module, we show how once trained our system can continuously improve by ingesting novel patterns.

1. Introduction

What makes humans capable of succeeding in a large variety of tasks is the capacity to learn from experience, recalling past events and generalizing to new ones. Learning to drive is a clear example of this ability. In recent years a lot of effort has been made to imitate this skill and to develop autonomous vehicles that are able to safely drive among other agents, either autonomous or driven by humans. Whereas remarkable progress has been made for automotive [2, 8, 38], current approaches still lack the ability to explicitly remember specific instances from experience when trying to infer possible future states of surrounding agents. This is particularly important for predicting future locations of moving agents, so to take appropriate decisions and avoid collisions or potentially dangerous situations. Predicting future trajectories of such agents is intrin-

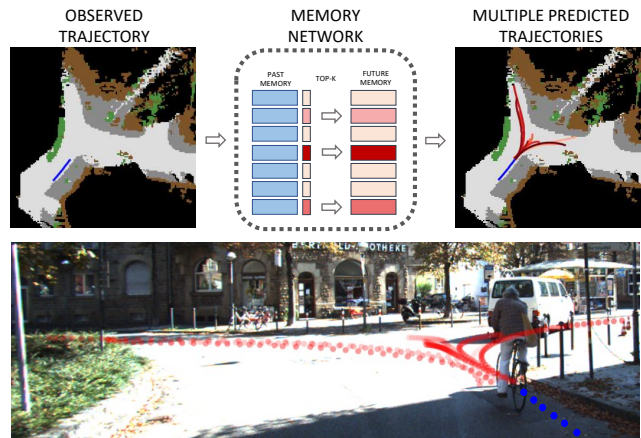


Figure 1. MANTRA addresses multimodal trajectory prediction. We obtain multiple future predictions given an observed past relying on a Memory Augmented Neural Network.

sically multimodal: vehicle dynamics give rise to a set of similarly likely outcomes for an external observer (Fig. 1).

While humans can address this task by implicit learning, i.e. exploiting procedural memory (knowing how to do things) from similar scenarios of previous experience, without explicit and conscious awareness, for machines this task has proven to be extremely hard. Common machine learning models, such as Recurrent Neural Networks, fail to address it. They are capable to store past information into an internal state, updated at every time step, and make predictions based on long term patterns. But in such networks, memory is a single hidden representation and is only addressable as a whole. State to state transition is unstructured and global. Instead, an element-wise addressable memory would be useful to selectively access only relevant pieces of information. This would allow to peak into likely futures to guide predictions.

In this paper we present MANTRA: Memory Augmented Neural TRAjectory predictor. MANTRA is a novel approach implementing a persistent Memory Augmented Neural Network (MANN) for vehicle trajectory prediction. In our model, an external associative memory is trained to write pairs of past and future trajectories and keep in mem-

ory only the most meaningful and non-redundant samples.

The model incrementally creates a knowledge base that is used as experience to perform meaningful predictions. This mimics the way in which implicit human memory works. Since the knowledge base is built from trajectory samples, it can also include instances observed while the system is running, after it has been trained. In this way the system gains experience online increasing its accuracy and capability to generalise at no training cost.

To memorize samples, past and future trajectories are stored in the memory in an encoded form, separately. In fact, this permits to use the encoding of an observed trajectory as a memory key to read an encoded future and decode them jointly to generate a prediction. Therefore, the actual coordinates are obtained decoding a future read from memory, conditioning it with the observed past. In this way, the output is not a simple copy of previously seen examples, but is instead a newly generated trajectory obtained both from the system experience (i.e. its memory) and the instance observed so far. By reading multiple futures from memory, diverse meaningful predictions can be obtained. The main contributions of this paper are the following:

- We propose a novel architecture for multiple trajectory prediction based on Memory Augmented Neural Networks. To the best of our knowledge we are the first to adopt MANNs for trajectory prediction.
- Our formulation, exploiting an encoder-decoder pipeline augmented with an associative memory, is easier to inspect and provides naturally multimodal predictions, obtaining state-of-the-art results on three traffic datasets.
- Our model is able to improve incrementally, after it has been trained, when observing new examples online. This trait is important for industrial automotive applications and is currently lacking in other state of the art predictors.

2. Related Work

Trajectory Prediction Significant effort has been made in the past years regarding trajectory prediction. Several researchers have focused on trajectories of pedestrians [1, 13, 14, 27, 31], either regarded as individuals or crowds, also exploiting social behaviors and interactivity between individuals [1, 13, 14, 21, 27]. While relevant for pedestrians, social behaviors are much less relevant for vehicles [20]. In this context, focus shifts instead on the observation of motion of the individual agents (their past trajectory) and the understanding of the surrounding environment [20, 33]. Traffic dynamics likely reduce to simpler scenarios where movement is limited and constrained by the environment. A notable exception is estimating lane changes on highways [10, 18]. A few efforts have been made to understand and predict vehicle trajectories in urban scenarios [20, 23, 33, 41]. Among them, DESIRE [20]

uses a Variational Autoencoder for estimating a distribution from which future trajectories can be sampled. The method though is not able to generate confidence scores to provide a ranked set of trajectories. A large number of predictions is needed to cover all the search space and Inverse Optimal Control is then used to extract a final ranked subset. INFER [33] instead exploits a fully convolutional model that takes into account intermediate semantic representations and generates multimodal heatmaps of possible future locations, then looking for peaks of the distribution.

In our work we address prediction of multiple vehicle trajectories in urban scenes. Examples of contexts where such multiple predictions may be necessary are roundabouts and crossroads where vehicles might take different equally possible paths. Differently from DESIRE [20] our approach is able to directly estimate a small set of ranked trajectories which already exhibit sufficient diversity to cover multiple futures. Differently from INFER [33] we directly work with coordinates instead of heatmaps, providing a better spatial resolution and more precise predictions. Differently from both DESIRE and INFER, we train a Memory Augmented Neural Network model to generate multimodal trajectories, which to the best of our knowledge has never been used with this purpose. The usage of MANNs has two main advantages: (i) multiple futures can be read from memory for a given trajectory observation, making the model capable to predict multiple outcomes, complying to the multimodal nature of the problem; (ii) by retrieving a likely future from memory we can rely on an oracle that suggests what is going to happen in the near future.

A conceptually similar research direction to ours is the one of intention-based methods [4, 7, 30]. Here, some anchor information (such as trajectories, actions or locations) are predefined and then used to guide predictions after estimating a probability distribution over each candidate. In [30], predictions are conditioned by the state of a robot agent, for which a goal is given or estimated. The authors of [7] propose a model for intersections that generates a likelihood over 5 fixed map zones, entailing different motion patterns. In [4], anchor trajectories are created with k-means and random sampling over training data. To some extent, our memory entries can be interpreted as anchors encoding futures instead of intentions. However, we do not choose a reference agent to condition predictions or restrict the applicability to constrained scenarios.

In order to obtain meaningful predictions we also take context into account and its physical constraints. According to this, the set of trajectory proposals obtained by the MANN is refined by integrating knowledge of the surrounding environment using semantic maps. Finally, differently for prior work, our trajectory prediction model is also capable of growing online, improving incrementally its performance from new observations after it has been trained.

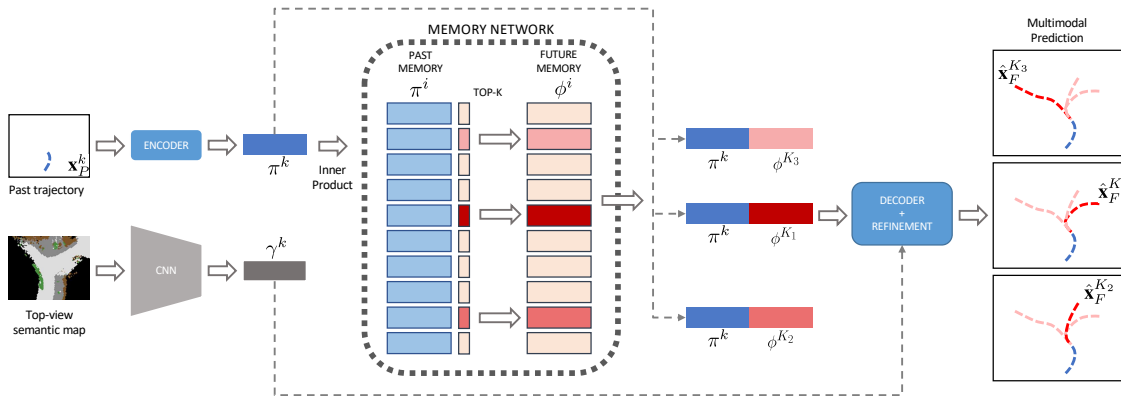


Figure 2. Architecture of MANTRA. The encoding of an observed past trajectory is used as key to read likely future encodings from memory. A multimodal prediction is obtained by decoding each future encoding, conditioned by the observed past. The surrounding context is processed by a CNN and fed to the Refinement Module to adjust predictions.

Memory Networks Neural networks with memory capabilities have been introduced to solve several machine learning problems which require to model a temporal dimension. The most common models are Recurrent Neural Networks (RNN) and their variants such as Long-Short Term Memories (LSTM) [15] and Gated Recurrent Units (GRU) [6]. However, in these models memory is a single hidden state vector that encodes all the temporal information. So memory is addressable as a whole and they lack the ability to address individual elements of knowledge, necessary to apply algorithmic manipulation and rapid inference. Moreover, state to state transition is unstructured and global. Being the state updated at each time-step, eventually it fails to model very long term dependencies. Finally, the number of parameters is tied to the size of the hidden state. So adding knowledge from the external environment necessarily implies increasing the size of the state. These characteristics prevent to use these models to effectively solve classes of problems like the one we address in this paper.

Recent works [3, 12, 16, 19, 22, 28, 29, 32, 34, 35, 37, 39, 40] have proposed Memory Augmented Neural Networks, or simply Memory Networks, to overcome the limitations of RNNs. The principal characteristic of this model is the usage of a controller network with an external element-wise addressable memory. This is used to store explicit information and access selectively relevant items. The memory controller is trained to dynamically managed memory content optimizing predictions. Differently from RNNs, state to state transitions are obtained through read/write operations and a set of independent states is maintained. An important consideration is that in Memory Networks the number of parameters is not tied to the size of the memory, i.e. increasing the memory slots will not increase the number of parameters.

While introduced recently, a number of applications of this model have already appeared in literature. The first em-

bodiment of a Memory Network was proposed in Neural Turing Machines (NTM) [12] to perform algorithmic tasks, such as sorting or copying, which require sequential manipulation steps. Thanks to a fully differentiable controller, the model interacts with the memory through read/write operations. The architecture was later extended to perform one-shot learning in [32]. Differently from NTM they trained the MANN to implement a Least Recently Used memory access strategy to write into rarely used locations.

In [37] MANNs have been proved to be able to effectively address Question Answering tasks, where the model has to answer questions related to a series of sentences. In [34] the same problem is solved with an End-to-End Memory Network with attention weights to shift importance from one sentence to another. Recent approaches have proposed a MANN to address the more complex problem of Visual Question Answering [19, 22], training the MANN to learn uncommon question-answer pairs. Online learning has also been tackled using Memory Networks. Rebuffi *et al.* [29] learn a classifier adding classes incrementally. MANNs for object tracking have been proposed where the model is trained to memorize templates, which are updated as the object is tracked [40].

All these MANNs rely on episodic memories. The system learns to write and read from memory but the stored data is limited only to the current set of observations (such as a list of numbers to be sorted in [12] or a collection of sentences for question answering in [37]). Differently from prior work, we build a MANN with a memory that is not episodic. Instead, it acts like a persistent memory which stores an experience of relevant data to perform accurate predictions for any observation and not just for a restricted episode or set of samples. The rationale behind this approach is that instead of solving simple algorithmic tasks as a Neural Turing Machine, we learn how to create a pool of samples to be used for future trajectory predictions.

The proposed model learns to store in memory only what is strictly needed to perform accurate predictions. Our usage of MANN is close to [26], but differs substantially. While they exploit the decoupling of embeddings to better fit data, we leverage the disjoint representation to create multiple outputs from a single input, leading to a fully multimodal predictive capability of the overall system.

3. Model

We formulate the task of vehicle trajectory prediction as the problem of estimating $P(\hat{\mathbf{x}}_F|\mathbf{x}_P, \mathbf{c})$, where $\hat{\mathbf{x}}_F$ is the predicted future trajectory, \mathbf{x}_P is the observed trajectory (or *past*) and \mathbf{c} is a representation of the context (e.g. roads, sidewalks). We consider vehicle trajectories as a sequence of 2-dimensional spatial coordinates. The *past* \mathbf{x}_P is given by its positions observed up to some reference point identified as *present*. Similarly, the *future* \mathbf{x}_F is the sequence of positions in which it will find itself at the next time steps.

3.1. Memory Based Trajectory Prediction

Given a sample trajectory $\mathbf{x}^i = [\mathbf{x}_P^i, \mathbf{x}_F^i]$, let $\pi^i = \Pi(\mathbf{x}_P^i)$ and $\phi^i = \Phi(\mathbf{x}_F^i)$ be two encoding functions that map the 2D coordinates of past and future trajectories into two separate latent representations. Similarly, let $\Psi(\pi^i, \phi^i)$ be a function that decodes a pair of past-future encodings into the coordinates of the future sub-trajectory \mathbf{x}_F^i , as shown in Fig. 2.

We define $M = \{\pi^i, \phi^i\}$ as an associative key-value memory containing $|M|$ pairs of past-future encodings. When a new trajectory \mathbf{x}_P^k is observed, its encoding π^k is used as key to retrieve meaningful samples from memory. Note that observed trajectories are all considered to be *past* trajectories, since the future counterpart is yet to be observed and is what we want to predict. The memory addressing mechanism is implemented as a cosine distance between past encodings, which produces similarity scores $\{s_i\}$ over all memory locations:

$$s_i = \frac{\pi^k \pi^i}{\|\pi^k\| \|\pi^i\|} \quad i = 0, \dots, |M| \quad (1)$$

According to the similarity scores, the future encodings of the top-K elements ϕ^j are separately combined with the encoding of the observed past π^k . The novel pairs of encodings are transformed into 2D coordinates using the decoding function Ψ : $\hat{\mathbf{x}}_F^j = \Psi(\pi^k, \phi^j)$, with $j = 1, \dots, K$. Note that π^k is fixed while ϕ^j varies depending on the sample read from memory. Future encodings ϕ^j act as an oracle which suggests possible outcomes based on the past observation. This strategy allows the model to look ahead into likely futures in order to predict the correct one. Since multiple ϕ^j can be used independently, we can decode multiple futures and obtain a multimodal prediction in case of uncertainty (e.g. a bifurcation in the road).

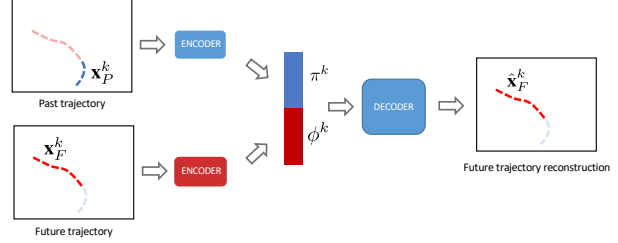


Figure 3. Representation learning: past and future trajectories are encoded separately; a decoder reconstructs future trajectory only.

3.2. Feature Representation Learning

The encoding-decoding functions Π, Φ, Ψ are trained jointly as an autoencoder, as shown in Fig. 3. The encoders learn to map past and future points into a meaningful representation and the decoder learns to reproduce the future. Instead of using just the future as input, we condition the reconstruction process also with an encoding of the past. This is useful for two aspects. First, we are able to train two different encoders for past and future. The two encoders are used to obtain separate representations for both keys (past) and values (future) in memory. Second, we obtain reconstructions of the future that is compatible with the past. This is of crucial importance for prediction since at test time we synthesize trajectory encodings by combining past and future parts taken from different examples. This also allows to generate trajectories that differ from the ones in memory and are not just a simple copy of already observed samples.

3.3. Memory controller

Traditional Memory Augmented Neural Networks [12, 34, 37] are designed to observe collections of data, usually referred to as episodes. The models are equipped with a working memory to store relevant information about the episode in order to generate a meaningful output for the episode. Yet memory is cleared for each episode and what is trained is the controller that decides what to read/write. The supervision for training stems from the cost function at the end of the episode, tracing gradients back to the controller.

As in standard memories, we train a controller to emit a write probability $P(w)$ every time that a sample is observed but, differently from these approaches, we train it to build a compact and expressive permanent memory. Training such a controller might result challenging since $P(w)$ does not depend only on the intrinsic importance of the observed sample but also on the current state of the memory. To solve this issue, we do not rely on the prediction loss for supervision. We instead feed the reconstruction error e to the controller, which decides if the network was sufficiently close to the ground truth. To enforce this behavior we define the controller loss \mathcal{L}_c as:

$$\mathcal{L}_c = e \cdot (1 - P(w)) + (1 - e) \cdot P(w) \quad (2)$$

where e is assumed to have values in $[0, 1]$. When the error is low, i.e. $e \rightarrow 0$, then

$$\mathcal{L}_c \approx P(w) \quad (3)$$

therefore the write probability is minimized.

Conversely, when $e \rightarrow 1$, then

$$\mathcal{L}_c \approx 1 - P(w) \quad (4)$$

and the controller maximizes the write probability.

What the controller is learning is an adaptive threshold on the reconstruction error, that allows to store in memory only what is useful to predict accurately, limiting redundancy. If the model exhibits a large prediction error, the controller writes the current sample with its ground truth future encoding in memory. When this happens, it indicates that the memory lacks samples to accurately reconstruct the future. Hence, by writing the sample in memory, the model will improve its prediction capabilities.

To satisfy the assumption of a bounded error function with values in $[0, 1]$ for the controller loss of Eq. 2, we introduce an adaptive miss rate error function with a threshold depending on the timestep:

$$e = 1 - \frac{1}{N} \sum_{i=1}^N \mathbb{1}_i(\hat{\mathbf{x}}_F, \mathbf{x}_F) \quad (5)$$

where $\mathbb{1}_i(\hat{\mathbf{x}}_F, \mathbf{x}_F)$ is an indicator function equal to 1 if the i -th point of the prediction $\hat{\mathbf{x}}_F$ lays within a threshold th from the ground truth and 0 otherwise. We use a different threshold for each timestep, allowing a given uncertainty for the farthest point (4 seconds) and linearly decreasing towards 0. In our experiments we use $th_{4s} = 2m$.

3.4. Iterative Refinement Module

To ensure compatibility with the environment, we refine predictions with an iterative procedure. Similarly to DESIRE [20], we adopt a feature pooling strategy: first, a CNN extracts a feature map γ^k from the context \mathbf{c} ; then, predictions are overlapped with the feature map and, for each time step coordinates, we extract the correspondent feature values (one per channel); finally, the resulting vector is fed to a GRU and a fully connected that output trajectory offsets.

The CNN is: $8 \times (k3, s2, p1)$; $16 \times (k3, s1, p1)$, where k is kernel size, s stride, p padding. Both layers have Batch-Norm and ReLU. The GRU has a hidden state size of 48. We do 4 iterations and we observed that increasing them does not introduce substantial changes.

3.5. Training

We train our model to observe 2 seconds trajectories and predict 4 seconds in the future. To achieve translation and rotation invariance, each trajectory is normalized by shifting

the present in the origin and rotating the trajectory in order to make it tangent with the Y-axis in the origin. In this way all futures start from $(0, 0)$ in an upward direction.

First, a pretraining of both the encoders and the decoder is done jointly as an autoencoder. To do so, we feed pairs of past and future trajectories belonging to the same samples, reconstructing only the future coordinates. We then train the memory controller, exploiting the learned past encoder and future decoder and resetting memory after every epoch. As controller we use a linear layer with sigmoid activation. The trained controller allows the memory to be filled with useful and non-redundant training samples by iterating over the training set and measuring their reconstruction error. While in principle the order in which samples are presented to the memory for writing may result in different final content, in our experiments we found that this does not affect the final prediction result. As a last step, we jointly train the refinement module and finetune the decoder. Here we feed the decoder with past and future encodings belonging to different samples, since the future is read from memory.

The two encoders and the decoder are implemented as Gated Recurrent Units with a 48-dimensional hidden state for each encoder and 96-dimensional for the decoder. The GRU in the refinement module is initialized with the past embedding and takes as input the predicted coordinates. This provides the module with complete information about the whole trajectory. We optimize \mathcal{L}_c defined in Eq. 2 to train the controller and a Mean Squared Error loss for decoder and refinement. All components are trained with the Adam optimizer using a learning rate of 0.0001.

4. Experiments

4.1. Datasets

KITTI [11] The dataset includes many annotations such as Velodyne LiDAR 3D scans, object bounding boxes and tracks, calibration, depth and IMU. Not all data is present for every video so we used the ones categorized as *KITTI Raw Data*, following the split of DESIRE [20]. Although the split is known, how to divide trajectories in data chunks is not. To obtain samples we collect 6 seconds chunks (2 seconds for past and 4 for future) in a sliding window fashion from all trajectories in the dataset, including the ego-vehicle. We obtain 8613 top-view trajectories for training and 2907 for testing. Note that these numbers are different from the original DESIRE split since they claim to gather 2509 trajectories in total. To favor reproducibility and future comparison we will publicly release our version of the dataset upon publication. Since top-view maps are not provided by KITTI, we project semantic labels of static categories obtained with DeepLab-v3+ [5] from all frames in a common top-view map using the Velodyne 3D point cloud and IMU. The resulting maps have a spatial resolution of 0.5 meters, and will be released along with the trajectories.

Method	ADE				FDE			
	1s	2s	3s	4s	1s	2s	3s	4s
Kalman	0.51	1.14	1.99	3.03	0.97	2.54	4.71	7.41
Linear	0.20	0.49	0.96	1.64	0.40	1.18	2.56	4.73
MLP	0.20	0.49	0.93	1.53	0.40	1.17	2.39	4.12
MANTRA (top 1)	0.24	0.57	1.08	1.78	0.44	1.34	2.79	4.83
MANTRA (top 5)	0.17	0.36	0.61	0.94	0.30	0.75	1.43	2.48
MANTRA (top 10)	0.16	0.30	0.48	0.73	0.26	0.59	1.07	1.88
MANTRA (top 20)	0.16	0.27	0.40	0.59	0.25	0.49	0.83	1.49
DESIRE (top 1) [20]	-	-	-	-	0.51	1.44	2.76	4.45
DESIRE (top 5) [20]	-	-	-	-	0.28	0.67	1.22	2.06
DESIRE (top 20) [20]	-	-	-	-	-	-	-	2.04

Table 1. Results on the KITTI dataset. Results obtained by DESIRE are given as reference even if not comparable, due to the data collection process.

Another smaller version of KITTI for trajectory prediction has been recently proposed by [33] and is publicly available. The authors propose 5 different train/test splits and average results over all runs, so we follow this evaluation protocol. We report experiments on both variants of KITTI. In the following we will refer to KITTI as our split obtained following DESIRE, unless stated otherwise.

Oxford RobotCar [25] & Cityscapes [9] The two datasets RobotCar and Cityscapes have been adapted for trajectory prediction in [33] to show zero-shot transfer capabilities on different domains. Of particular interest is the ability to transfer to RobotCar since the sequences are acquired in the UK where cars drive on the left-side of the road. RobotCar has 6 seconds trajectories divided into 2 seconds for past and 4 for future. Cityscapes instead has shorter videos and predictions are made only up to one second in the future, as done in [33].

4.2. Evaluation metrics and Baselines

We report results in two common metrics for vehicle trajectory prediction: *Average Displacement Error* (ADE) and *Final Displacement Error* (FDE), where ADE is the average L2 error between all future timesteps and FDE (sometimes referred to as Horizon error) is the error at a given timestep. As in [20, 33] we take the best out of K predictions to account for the intrinsic multimodality of the task. We compare our approach with several baselines: a linear coordinate regressor (*Linear*); a Multi-Layer Perceptron with two layers trained as a coordinate regressor (*MPL*); a Kalman filter [17], with a constant speed model used to propagate the estimate without incorporating measures (*Kalman*). We implemented and tested the baselines on the KITTI dataset to show comparable results. When available we also report existing baselines from the literature.

4.3. Results

Table 1 shows the results on the KITTI dataset. Simply propagating the trajectory with a Kalman filter proves to be insufficient to accurately predict future positions, especially over long time spans, with a FDE@4s higher than 7m.

Method	ADE				FDE			
	1s	2s	3s	4s	1s	2s	3s	4s
Kalman	0.33	0.54	0.93	1.4	0.46	1.18	2.18	3.32
Linear	0.31	0.56	0.89	1.28	0.47	1.13	1.94	2.87
MLP	0.30	0.54	0.88	1.28	0.46	1.12	1.94	2.88
RNN Enc-Dec [36]	0.68	1.94	3.20	4.46	-	-	-	-
Markov [33]	0.70	1.41	2.12	2.99	-	-	-	-
Conv-LSTM (top 5) [33]	0.76	1.23	1.60	1.96	-	-	-	-
INFER (top 1) [33]	0.75	0.95	1.13	1.42	1.01	1.26	1.76	2.67
INFER (top 5) [33]	0.56	0.75	0.93	1.22	0.81	1.08	1.55	2.46
MANTRA (top 1)	0.37	0.67	1.07	1.55	0.60	1.33	2.32	3.50
MANTRA (top 5)	0.33	0.48	0.66	0.90	0.45	0.78	1.22	2.03
MANTRA (top 10)	0.31	0.43	0.57	0.78	0.43	0.67	1.04	1.78
MANTRA (top 20)	0.29	0.41	0.55	0.74	0.41	0.64	1.00	1.68

Table 2. Results on the KITTI dataset (INFER split).

Learning based baselines all perform better than Kalman filter, with the Multi-Layer Perceptron performing slightly better than the linear regressor.

Models that generate a single prediction fail to address the multimodality of the task, since they are trained to lower the error with a single output even when there might be multiple equally likely desired outcomes. What may happen is that in front of a bifurcation, the model predicts an average of the two possible trajectories, trying to satisfy both scenarios. Examples of this behavior are shown in Fig. 4. Each prediction of MANTRA instead follows a specific path, ignoring the others. This leads to high errors on some examples when generating only one future, since the model may decide to follow a different likely path. On the other hand as soon as we generate K multiple predictions, the top- K error drastically decreases since we are able to cover diverse future paths. We also report results from DESIRE [20] varying K . Even though these results are not directly comparable as explained in Section 4.1, it is interesting to observe how DESIRE quickly saturates when increasing K , while our method keeps lowering the error significantly. This suggests that MANTRA samples a higher diversity of futures both at a coarse level (i.e. taking one road or another) and at a fine level (i.e. different behaviors on the same road). Some qualitative results on KITTI are shown in Fig 4, comparing them with the baselines.

Additionally, we evaluate MANTRA on the KITTI split proposed in [33], as shown in Table 2. Here we also report some available baselines from the state of the art, both for single and multimodal predictions. With $K = 1$ our method performs better or on par with INFER [33] at low timesteps, yet losing some precision at 4s. Increasing K instead we are able to largely outperform INFER over all timesteps.

Following [33], we showcase the ability of our model to zero-shot transfer to other datasets. On Oxford RobotCar (Tab. 3) MANTRA is still able to provide satisfactory results, consistently outperforming INFER across timesteps for multimodal predictions. Analogously, on Cityscapes (Tab. 4) the model obtains a lower error than the other methods. Here we report only errors at 1s in the future, which is the maximum length of the trajectories in the dataset.

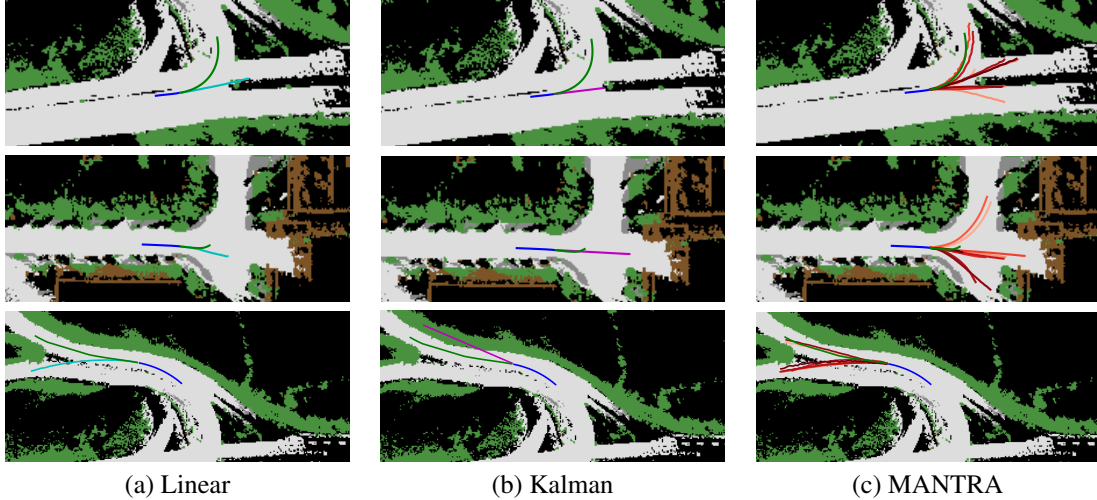


Figure 4. MANTRA compared to Linear regression (a) and Kalman filter (b). Methods (a),(b) lack multi-modal capability. Past trajectories are depicted in blue, ground truth in green and future predictions are cyan (a), purple (b) and red (c). In (c) highly ranked are darker.

Method	ADE				FDE			
	1s	2s	3s	4s	1s	2s	3s	4s
INFER (top 1) [33]	1.06	1.35	1.48	1.68	1.31	1.71	1.70	2.56
INFER (top 5) [33]	0.85	1.14	1.29	1.50	1.18	1.58	1.58	2.41
MANTRA (top 1)	0.55	0.77	1.01	1.30	0.60	1.15	1.82	2.63
MANTRA (top 5)	0.55	0.68	0.82	1.03	0.58	0.88	1.37	2.07
MANTRA (top 10)	0.44	0.56	0.72	0.94	0.48	0.73	1.33	1.98
MANTRA (top 20)	0.31	0.43	0.59	0.83	0.35	0.61	1.24	1.96

Table 3. Results on the Oxford RobotCar dataset.

Method	ADE	FDE
Conv-LSTM (top 1) [33]	1.50	-
Conv-LSTM (top 3) [33]	1.36	-
Conv-LSTM (top 5) [33]	1.28	-
INFER (top 1) [33]	1.11	1.59
INFER (top 3) [33]	0.99	1.45
INFER (top 5) [33]	0.91	1.38
MANTRA (top 1)	0.81	1.42
MANTRA (top 3)	0.66	1.15
MANTRA (top 5)	0.60	1.00
MANTRA (top 10)	0.54	0.86
MANTRA (top 20)	0.49	0.79

Table 4. Results on the Cityscapes dataset at 1s in the future.

4.4. Incremental setting

Differently from prior work on trajectory prediction, MANTRA is able to improve its capabilities online, i.e. observing other agents' behaviors while driving. We simulate an online scenario on KITTI, iteratively removing a small set of 50 trajectories from the test set, presenting them to the memory controller. The controller incorporates novel patterns according to $P(w)$. At each iteration we test the predictor on the remaining test set. In Fig. 5 memory growth and test error are shown for MANTRA with $K=5$ multiple futures. Similar behaviors can be observed varying K . Interestingly, the memory size slowly grows while the error keeps decreasing. Note that the memory stores only the 16% of the newly seen examples. To cope with the error variance increase when the remaining set of samples decreases in size we average results over 100 runs.

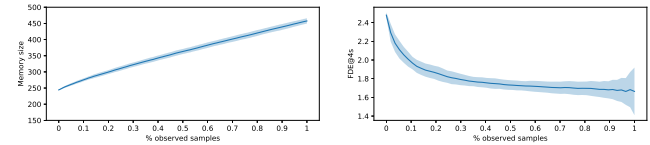


Figure 5. Online setting. Mean and variance of memory growth (left) and error rate (right) averaged over 100 runs, increasing the observed samples.

5. Model Analysis

In the following we perform ablation studies aimed at highlighting the importance of components in our model. We thoroughly investigate how the model organizes memory by checking what gets written and how it gets decoded.

Ablation Studies We investigate modifications of MANTRA reporting results in Tab. 5 on KITTI. We test the following: (i) without refinement; (ii) without decoder, i.e. reading from memory using encodings but just copying the correspondent future coordinates; (iii) without rotation invariance, i.e. using trajectories with random rotations; (iv) without memory controller, i.e. adding all training samples in memory; (v) without encoder-decoder, i.e. a Nearest Neighbor between past trajectory coordinates copying the future of the closest sample in coordinate space. On the one hand, when the memory is filled with all training samples instead of selecting them with a controller, the error drastically increases; on the other hand even worse results are obtained when the samples are not encoded and decoded with the recurrent GRU layers. Even removing just the decoder lowers the precision of predictions considerably. This should not come as a surprise, since the decoder has the important role of adapting the suggested future from memory to the current sample, making it coherent with its past. Surprisingly enough instead, the

Method	ADE				FDE				Memory Size
	1s	2s	3s	4s	1s	2s	3s	4s	
MANTRA (top 5)	0.17	0.36	0.61	0.94	0.30	0.75	1.43	2.48	190 (2.2 %)
MANTRA w/o ref.	0.18	0.39	0.67	1.04	0.33	0.85	1.59	2.65	190 (2.2 %)
MANTRA w/o dec.	0.25	0.46	0.76	1.18	0.42	0.91	1.75	3.12	190 (2.2 %)
MANTRA w/o rot. inv.	0.25	0.51	0.88	1.38	0.45	1.09	2.10	3.58	2170 (25.2 %)
MANTRA w/o ctrl.	0.20	0.45	0.82	1.34	0.37	1.02	2.07	3.64	8613 (100 %)
MANTRA w/o enc-dec.	0.24	0.58	1.08	1.75	0.47	1.36	2.74	4.68	8613 (100 %)

Table 5. Ablation study of the full method against variants without specific components: decoder, refinement, rotation invariance, trained controller, encoder-decoder. Errors are at K=5. Memory size is shown as number of samples and % of the training set.

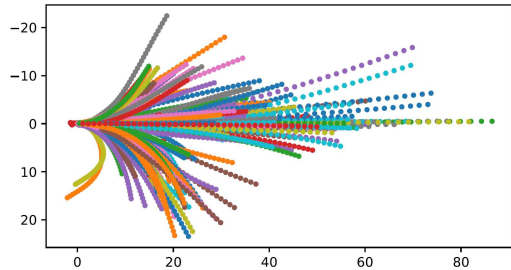


Figure 6. Decoded trajectories from memory.

refinement module does not play a very important role in the reconstruction, suggesting that the originally generated trajectories are already precise. Rotation invariance proves to be very relevant in moderating memory size and improving accuracy. By adding rotation invariance to training we lower the memory size from the 25.2% to the 2.2% of the observed training set.

Memory inspection To understand what the model is learning, we inspect what the controller stores in memory. We take each sample and plot its decoded future to depict a snapshot of every sample in memory. Fig. 6 shows all samples from a memory filled for K=5 predictions.

In Fig. 7 we plot T-SNE projections [24] of past and future encodings in memory, as points. On the left we plot past embeddings, while on the right we report future embeddings. For each projected sample we show future trajectories generated by the decoder, displayed starting from T-SNE points. All trajectories in the image have an upward trend due to the rotation invariance we introduce for storing samples. Similar trajectories are clustered together, indicating that the encoders are learning a manifold where similar patterns are close. Observing the T-SNE of past encodings, the multimodal nature of the problem emerges. In fact, the space appears to be organized mostly by trajectory speed and for each point several possible future directions are present. When trajectories have lower speed, futures are free to span over many possible directions, while when trajectories have higher speed, the futures vary more in length rather than curvature.

Decoder Analysis We inspect the behavior of the decoder and the influence that different pasts have on future reconstructions. Encoder and decoder are jointly trained, but differently from standard autoencoders, only part of the input

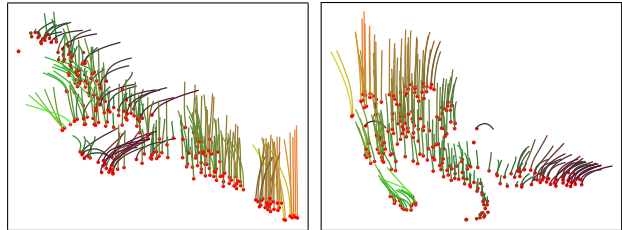


Figure 7. T-SNE representations of past (left) and future (right) encodings stored in memory. Each point in the embedding space is shown along with the decoded trajectory. Trajectories are color coded by orientation (green tones) and speed (red tones).

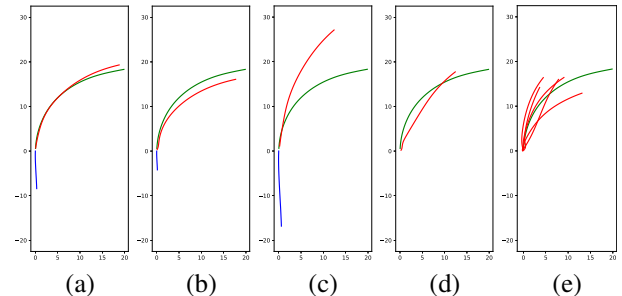


Figure 8. Influence of past in the decoder. (a) observed past; (b) slower past; (c) faster past; (d) past embedding zeroed; (e) multiple randomized past embeddings. Blue: past trajectory used for decoding. Red: future reconstruction. Green: original future.

is reconstructed, i.e. the future. The past has the important role of conditioning the reconstruction so that we can generalize to unseen examples. In Fig. 8 we show several reconstructions of the same future, changing only the past encoding and keeping fixed the future one. The reconstructions of the original past yields a precise reconstruction. By changing the past by shortening it or stretching it, i.e. changing the velocity, the reconstruction gets accelerated or decelerated, affecting curvature. As a control experiment we also use a vector of zeros or a random embedding. In both cases the generated trajectories are very imprecise but still follow approximately the original trend. These tests justify using the decoder feeding a combination of encodings belonging to different samples, as we do at test time. In fact the generated trajectories are new compared to the samples in memory and they adapt to the current observation.

6. Conclusions

We propose MANTRA, the first Memory Augmented Neural TRAjectory prediction framework. Our method, based on an associative memory, can natively grasp the inherently multi-modal nature of the future trajectory prediction problem, yielding state of the art results on three traffic datasets. Moreover, we show that the memory is able to ingest novel samples lowering the error on unseen data.

Acknowledgments We thank NVIDIA for donating a Titan Xp GPU. This work is partially funded by IMRA Europe S.A.S.

References

- [1] Alexandre Alahi, Kratharth Goel, Vignesh Ramanathan, Alexandre Robicquet, Li Fei-Fei, and Silvio Savarese. Social lstm: Human trajectory prediction in crowded spaces. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 961–971, 2016.
- [2] Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Prasoon Goyal, Lawrence D Jackel, Mathew Monfort, Urs Muller, Jiakai Zhang, et al. End to end learning for self-driving cars. *arXiv preprint arXiv:1604.07316*, 2016.
- [3] Qi Cai, Yingwei Pan, Ting Yao, Chenggang Yan, and Tao Mei. Memory matching networks for one-shot image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4080–4088, 2018.
- [4] Yuning Chai, Benjamin Sapp, Mayank Bansal, and Dragomir Anguelov. Multipath: Multiple probabilistic anchor trajectory hypotheses for behavior prediction. *arXiv preprint arXiv:1910.05449*, 2019.
- [5] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*, 2017.
- [6] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.
- [7] Chiho Choi, Abhishek Patil, and Srikanth Malla. Drogon: A causal reasoning framework for future trajectory forecast. *arXiv preprint arXiv:1908.00024*, 2019.
- [8] Felipe Codevilla, Matthias Müller, Antonio López, Vladlen Koltun, and Alexey Dosovitskiy. End-to-end driving via conditional imitation learning. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1–9. IEEE, 2018.
- [9] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3213–3223, 2016.
- [10] Nachiket Deo and Mohan M Trivedi. Multi-modal trajectory prediction of surrounding vehicles with maneuver based lstms. In *2018 IEEE Intelligent Vehicles Symposium (IV)*, pages 1179–1184. IEEE, 2018.
- [11] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 3354–3361. IEEE, 2012.
- [12] Alex Graves, Greg Wayne, and Ivo Danihelka. Neural Turing machines. *arXiv preprint arXiv:1410.5401*, 2014.
- [13] Agrim Gupta, Justin Johnson, Li Fei-Fei, Silvio Savarese, and Alexandre Alahi. Social gan: Socially acceptable trajectories with generative adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2255–2264, 2018.
- [14] Dirk Helbing and Peter Molnar. Social force model for pedestrian dynamics. *Physical review E*, 51(5):4282, 1995.
- [15] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [16] Łukasz Kaiser, Ofir Nachum, Aurko Roy, and Samy Bengio. Learning to remember rare events. *arXiv preprint arXiv:1703.03129*, 2017.
- [17] Rudolph Emil Kalman. A new approach to linear filtering and prediction problems. 1960.
- [18] Alex Kuefler, Jeremy Morton, Tim Wheeler, and Mykel Kochenderfer. Imitating driver behavior with generative adversarial networks. In *2017 IEEE Intelligent Vehicles Symposium (IV)*, pages 204–211. IEEE, 2017.
- [19] Ankit Kumar, Ozan Irsoy, Peter Ondruska, Mohit Iyyer, James Bradbury, Ishaan Gulrajani, Victor Zhong, Romain Paulus, and Richard Socher. Ask me anything: Dynamic memory networks for natural language processing. In *International conference on machine learning*, pages 1378–1387, 2016.
- [20] Namhoon Lee, Wongun Choi, Paul Vernaza, Christopher B Choy, Philip HS Torr, and Manmohan Chandraker. Desire: Distant future prediction in dynamic scenes with interacting agents. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 336–345, 2017.
- [21] Matteo Lisotto, Pasquale Coscia, and Lamberto Ballan. Social and scene-aware trajectory prediction in crowded spaces. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 0–0, 2019.
- [22] Chao Ma, Chunhua Shen, Anthony Dick, Qi Wu, Peng Wang, Anton van den Hengel, and Ian Reid. Visual question answering with memory-augmented networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6975–6984, 2018.
- [23] Yuexin Ma, Xinge Zhu, Sibozhang, Ruigang Yang, Wenping Wang, and Dinesh Manocha. Trafficpredict: Trajectory prediction for heterogeneous traffic-agents. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 6120–6127, 2019.
- [24] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(Nov):2579–2605, 2008.
- [25] Will Maddern, Geoffrey Pascoe, Chris Linegar, and Paul Newman. 1 year, 1000 km: The oxford robotcar dataset. *The International Journal of Robotics Research*, 36(1):3–15, 2017.
- [26] Alexander Miller, Adam Fisch, Jesse Dodge, Amir-Hossein Karimi, Antoine Bordes, and Jason Weston. Key-value memory networks for directly reading documents. *arXiv preprint arXiv:1606.03126*, 2016.
- [27] Stefano Pellegrini, Andreas Ess, Konrad Schindler, and Luc Van Gool. You’ll never walk alone: Modeling social behavior for multi-target tracking. In *2009 IEEE 12th International Conference on Computer Vision*, pages 261–268. IEEE, 2009.
- [28] Alexander Pritzel, Benigno Uribe, Sriram Srinivasan, Adria Puigdomenech Badia, Oriol Vinyals, Demis Hassabis, Daan Wierstra, and Charles Blundell. Neural episodic control. In *Proceedings of the 34th International Conference*

- on *Machine Learning-Volume 70*, pages 2827–2836. JMLR.org, 2017.
- [29] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl: Incremental classifier and representation learning. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 2001–2010, 2017.
 - [30] Nicholas Rhinehart, Rowan McAllister, Kris Kitani, and Sergey Levine. Precog: Prediction conditioned on goals in visual multi-agent settings. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2821–2830, 2019.
 - [31] Amir Sadeghian, Vineet Kosaraju, Ali Sadeghian, Noriaki Hirose, Hamid Rezaatofghi, and Silvio Savarese. Sophie: An attentive gan for predicting paths compliant to social and physical constraints. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1349–1358, 2019.
 - [32] Adam Santoro, Sergey Bartunov, Matthew Botvinick, Daan Wierstra, and Timothy Lillicrap. Meta-learning with memory-augmented neural networks. In *International conference on machine learning*, pages 1842–1850, 2016.
 - [33] Shashank Srikanth, Junaid Ahmed Ansari, Sarthak Sharma, et al. Infer: Intermediate representations for future prediction. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS 2019)*, 2019.
 - [34] Sainbayar Sukhbaatar, Jason Weston, Rob Fergus, et al. End-to-end memory networks. In *Advances in neural information processing systems*, pages 2440–2448, 2015.
 - [35] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks for one shot learning. In *Advances in neural information processing systems*, pages 3630–3638, 2016.
 - [36] Jaskaran Viridi. *Using deep learning to predict obstacle trajectories for collision avoidance in autonomous vehicles*. PhD thesis, UC San Diego, 2017.
 - [37] Jason Weston, Sumit Chopra, and Antoine Bordes. Memory networks. *arXiv preprint arXiv:1410.3916*, 2014.
 - [38] Huazhe Xu, Yang Gao, Fisher Yu, and Trevor Darrell. End-to-end learning of driving models from large-scale video datasets. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2174–2182, 2017.
 - [39] Zhongwen Xu, Linchao Zhu, and Yi Yang. Few-shot object recognition from machine-labeled web images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1164–1172, 2017.
 - [40] Tianyu Yang and Antoni B Chan. Learning dynamic memory networks for object tracking. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 152–167, 2018.
 - [41] Alex Zyner, Stewart Worrall, and Eduardo Nebot. Naturalistic driver intention and path prediction using recurrent neural networks. *IEEE Transactions on Intelligent Transportation Systems*, 2019.