

LTN: Long-Term Network for Long-Term Motion Prediction

YingQiao Wang

Abstract—Making accurate motion prediction of surrounding agents such as pedestrians and vehicles is a critical task when robots are trying to perform autonomous navigation tasks. Recent research on multi-modal trajectory prediction, including regression and classification approaches, perform very well at short-term prediction. However, when it comes to long-term prediction, most Long Short-Term Memory (LSTM) based models tend to diverge far away from the ground truth. Therefore, in this work, we present a two-stage framework for long-term trajectory prediction, which is named as Long-Term Network (LTN). Our Long-Term Network integrates both the regression and classification approaches. We first generate a set of proposed trajectories with our proposed distribution using a Conditional Variational Autoencoder (CVAE), and then classify them with binary labels, and output the trajectories with the highest score. We demonstrate our Long-Term Network’s performance with experiments on two real-world pedestrian datasets: ETH/UCY, Stanford Drone Dataset (SDD), and one challenging real-world driving forecasting dataset: nuScenes. The results show that our method outperforms multiple state-of-the-art approaches in long-term trajectory prediction in terms of accuracy.

Index Terms—Trajectory prediction, long short-term memory (LSTM), robots, autonomous vehicles

I. INTRODUCTION

ACCURATELY predicting the motions of surrounding agents such as pedestrians and vehicles are significant when mobile robots or autonomous vehicles are trying to perform navigation tasks. In traffic systems, the future behavior of each traffic participant is determined by multiple aspects, such as the movement of other traffic agents, the physical constraints, and the traffic rules [1]–[3]. Humans have the ability to navigate through a complex traffic scenario because they have the ability to reason about all the other people’s actions, and how the physical constraints in the traffic systems affect their movements. Therefore, for a robot navigating through a complex traffic system, we need to consider about all the movement of the other surrounding traffic agents and the physical constraints in the traffic system.

With the discovery of the vanilla LSTM model, the researchers started to use Long Short-Term Memory networks to produce a regression of the future trajectories of traffic agents. The LSTM is a model that processes the data sequentially, so it is suitable for predicting the trajectories which is also considered as sequential data. Starting from the Social LSTM [4] model, the researchers started to model the people’s social interaction. When predicting the future trajectories of traffic agents, they will store the knowledge about people, e.g. speed,

direction, motion pattern, and people’s social interaction in the hidden state [5]–[8].

Then, the map information is integrated by extracting the features of map using a Convolutional Neural Network (CNN) combined with the current LSTM model, which largely improves the prediction accuracy. The recent works start to compete with each other by using different structures on modeling social interaction, and introduces a LSTM based encoder-decoder structure. The model encodes the past trajectory of the traffic agent using the LSTM along with the nearby traffic agents, produces a regression for the future trajectory, and decodes this trajectory using the LSTM. The Trajectron++ [2], WWTG [7] (Where Will They Go), and Social-BiGAT [9], which are recent state-of-the-art models, outperform most of the popular LSTM model on future trajectory prediction in terms of accuracy. The model uses the traditional LSTM encoder-decoder structure, but it encodes the past trajectory and future trajectory into a latent space using the Conditional Variational Autoencoder (CVAE) [10], [11]. For prediction, it draws a latent variable from the latent space, decodes it as a regression using GRU [12] and past trajectory information. However, there is still space for improvement in terms of accuracy in long-term prediction, especially 3 to 4 seconds after the current observation.

In this work, we propose a two-stage framework called Long-Term Network (LTN) to improve the long-term trajectory prediction in terms of accuracy. In the first stage, the LTN uses a traditional LSTM-GRU encoder-decoder structure along with the CVAE [10] to produce a set of possible future trajectory proposals. In the second stage, LTN performs classification and refinement on the trajectories proposals, and outputs the proposal with the highest score as the final trajectory prediction result. The trajectory proposals are generated based on the surrounding traffic agents identified by the LTN, and the prior extracted map information, so that the model can identify the traversable spaces of our robot and identify the possible effects of surrounding traffic agents to make better proposals.

The contributions of this paper are summarized as follows:

- 1) We propose a newly modified GRU unit called Mogrifier GRU, based on the idea of the Mogrifier LSTM [13]. By our refinement on the hidden state, we improve the performance of the model in terms of long-term prediction accuracy by 10% just by replacing the regular GRU with our Mogrifier GRU.
- 2) We propose a two-stage approach, in which we combine the regression and classification methods and largely improve the performance on the long term trajectory prediction.
- 3) Our model achieve the state-of-the-art results on the widely used

Y. Wang was with the Department of Mathematics and Statistics, Colby College, Maine, ME, 04901 USA e-mail: ywang22@colby.edu.

pedestrian trajectory prediction datasets (ETH/UCY) [14] [15], Stanford Drone Dataset (SDD) [16] and one real-world driving forecasting dataset (nuScenes) [17].

II. RELATED WORK

A. Multi-Modal Trajectory Prediction

Many earlier works in human trajectory forecasting can be roughly divided into two categories: the classic methods and the deep learning approaches. The classic methods include the kinematic equation, or the statistical models like polynomial fitting and Gaussian mixture models. In most recent works, Recurrent Neural Network (RNN) and its variants such as Long Short-Term Memory (LSTM) or Gated Recurrent Unit (GRU) [12], and Convolutional Neural Network become the basis of most of recent models. Researchers utilize CNN to extract map features, and RNN or its variants, LSTM and GRU, to capture the social interaction between the traffic agents, and then regress the future trajectory [9], [18]–[21].

B. Deep Generative Models

Besides the classic methods and the CNN-RNN combined approaches, the generative approaches have emerged as another state-of-the-art approach in trajectory prediction problems. This approach shifts some researchers' focus on regressing a single future trajectory to producing a distribution of the future trajectory. With a full distribution of the future trajectory, the results can produce more possible trajectory proposals. To produce such distributions, most works use recurrent backbone architecture with a latent variable model, such as the Generative Adversarial Network (GAN) [22], and the Conditional Variational Autoencoder (CVAE) [10]. Currently, Trajectron++ [2], Social-BiGAT [9], WWTG [7] are two CVAE and GAN based models that outperform most state-of-the-art trajectory prediction models. Trajectron++ [2] and Social-BiGAT [9] are able to account for the social interactions between traffic agents and physical constraints in the scene.

C. Regression and Classification

Current models that produce full distributions mostly utilize the Gaussian Mixture Model, which outputs the local maximum of the distribution as the final trajectory prediction result. But empirically, by the qualitative analysis in most of the work, the output is not actually the closest trajectory produced by the full distribution to the ground truth. So the new approaches combining regression and classification appear, which generates a set of hypothesis trajectory proposals, and outputs the proposal with the highest score as the final trajectory prediction result. Trajectory Proposal Net (TPNet) [19] is another state-of-the-art that uses this regression and classification method for trajectory prediction, where the model does polynomial fitting between the starting point and the proposed end point of the traffic agent, while considering the social-interaction and traffic rules. At the end, the model performs classification on these proposed trajectories and outputs the proposal with the top scores.

III. PROBLEM FORMULATION

In this work, we select our robot as the center point in the scene, where we will determine its surrounding traffic agents and predict their future trajectories. During the time interval $[0, T_{obs}] \cup [T_{obs} + 1, T_{future}]$, we denote the number of traffic agents surrounding the robot at time $t \in [0, T_{obs}] \cup [T_{obs} + 1, T_{future}]$ as S . We denote the surrounding agents as $\{A_1, \dots, A_S\}$, and for each agent, we categorize it as pedestrians or vehicles. For simplicity, the vehicle category also includes agents like bicycles, motorcycles, and cars. The model takes a series of past positions in the time interval $[0, T_{obs}]$ of each agent $P_{A_i}^{obs} = \{P_{A_i}^0, \dots, P_{A_i}^{T_{obs}}\}$, and also a series of past positions of the robot $P_r^{obs} = \{P_r^0, \dots, P_r^{T_{obs}}\}$, where $i \in [0, S]$. For each agent's future trajectory in interval $[T_{obs} + 1, T_{future}]$, we denote it as $P_{A_i}^{future} = \{P_{A_i}^{T_{obs}+1}, \dots, P_{A_i}^{T_{future}}\}$, and for the robot's future trajectory, we express it as $P_r^{future} = \{P_r^{T_{obs}+1}, \dots, P_r^{T_{future}}\}$, where $i \in [1, S]$. We also incorporate the map information in the same way as Trajectron++, and encode the map information as $M = M_S^t$, where the S is the agent, and $t \in [0, T_{obs}] \cup [T_{obs} + 1, T_{future}]$.

IV. METHOD

A. Long-Term Network

To further improve the performance of current model in long-term trajectory prediction, we propose a two-stage framework called Long-Term Network (LPN). The framework is visualized in Figure 1.

B. Determining the Surrounding Traffic Agents

To determine the surrounding traffic agents, we first determine the number of agents in the scene, and denote it as S . We include the agents A_i ($i \in [1, S]$) that are close to the robot in ℓ_2 distance. Formally, the agent is selected if at time $t \in [0, T_{obs}] \cup [T_{obs} + 1, T_{future}]$, $\|P_{A_i}^t - P_r^t\|^2 \leq d$, where d is a hyperparameter indicating the maximum perception distance.

Since we are going to predict each agent A_i 's future trajectories, we perform the similar process to select the surrounding traffic agents of our selected agents A_i . The agents around A_i is again determined by the ℓ_2 distance. Formally, the agent around A_i is selected if at time $t \in [0, T_{obs}] \cup [T_{obs} + 1, T_{future}]$ and $x \in [0, S]_{\neq i}$, $\|P_{A_i}^t - P_x^t\|^2 \leq d$, where d again is the same hyperparameter that expresses the maximum perception distance.

C. Modeling the Agent History and The Social Interactions

To model the agent history, we primarily utilize the Mogrifier LSTM [13], which is a variant of the vanilla LSTM model. The Mogrifier LSTM has better performance in long-term performance than the vanilla version, as the experiment in the paper demonstrates. The Mogrifier LSTM utilizes the same LSTM module, but between each unit, the Mogrifier LSTM updates the input and previous hidden state with several rounds of mutual gating, which is called a mogrifying step. The Mogrifier LSTM can also be implemented based on Bi-directional LSTM. Since there are no public code for the Mogrifier LSTM, we implement the bi-directional Mogrifier

the complete history tensor V_i , using the GRU [12] units. We present a more powerful variant of GRU unit, where we adopt the similar idea in Mogrifier LSTM [13] to process the input and hidden state before each GRU unit. Suppose we have input x_{input} and previous hidden state h_{prev} , for a normal GRU, the current hidden state h_{cur} is calculated by:

$$h_{cur} = GRU(x_{input}, h_{prev}), \quad (2)$$

and the h_{cur} is calculated by:

$$\begin{aligned} r_{cur} &= \sigma(W_{ir}x_{input} + b_{ir} + W_{hr}h_{prev} + b_{hr}) \\ z_{cur} &= \sigma(W_{iz}x_{input} + b_{iz} + W_{hz}h_{prev} + b_{hz}) \\ n_{cur} &= \tanh(W_{in}x_{input} + b_{in} + r_{cur} * (W_{hn}h_{prev} + b_{hn})) \\ h_{cur} &= (1 - z_{cur}) * n_{cur} + z_{cur} * h_{prev}, \end{aligned} \quad (3)$$

where $r_{cur}, z_{cur}, n_{cur}$ are the reset, update, and new gates. σ is the sigmoid function, and $*$ is the Hadamard product, and all the W, b are the learnable weights matrices.

Our Mogrifier GRU works by performing mogrifying steps before the usual GRU computation step. Suppose we perform mogrifying steps i times, we have $\text{MogGRU}(x_{input}, h_{prev}) = \text{GRU}(x_{input}^i, h_{prev}^i)$. For each $a \in [1, i]$:

$$x_{input}^a = 1.5 \tanh(Q^a h_{prev}^{a-1}) \odot x_{input}^{a-2}, \quad (4)$$

$$h_{prev}^a = 1.5 \tanh(R^a x_{input}^{a-1}) \odot h_{prev}^{a-2}, \quad (5)$$

where for odd $a \in [1, i]$, the equation (1) is performed, and for even $a \in [1, i]$, the equation (2) is performed. The parameters in the equations are recommended as: $i \in [5, 6]$ and i has to be a integer. In our model, we choose $i = 6$, and $i = 0$ will recover the traditional GRU unit. The parameter Q, R in the mogrifying steps are the same as the Q, R in the parameters in the Mogrifier LSTM, which are randomly initialized matrices. We will have the experiment result of Mogrifier GRU in the experiment section to show it's performance comparing with the traditional GRU.

G. Trajectory Proposal and Classification Module

During training, when we get the full distribution of the future trajectory of agent A_i , we then produce the final trajectory proposal by sampling N numbers of the trajectories from the latent variable distribution, where the latent variable is determined by:

$$z = \arg \max_{z \in Z} p_{\theta}(z|V_i, M). \quad (6)$$

Then, in the classification module. We consider the methods used in TPNNet [19]. We assign each proposal with a binary class label, which is used to indicate whether it is a good trajectory or not. We use the average distance between all the sampled trajectory proposals and the ground truth proposal as the criterion for measuring the proposal's quality, which can be expressed as:

$$D = \frac{1}{N} \sum_{n=1}^N \|p_{A_i}^n - p_{A_i}^{n,prop}\|^2, \quad (7)$$

which is the average ℓ_2 distance between the n -th sampled proposal vector $p_{A_i}^{n,prop}$ and the ground truth trajectory $p_{A_i}^{n,gt}$. Then, a threshold γ is used, and for the proposed trajectories that have D lower than γ , we assign positive labels. For the proposed trajectories that have D larger than γ , we assign negative labels.

H. Objective Function

In our model, there are two loss functions to be minimized, One is the regression loss, and the other one is the classification loss.

Regression Loss: We adopt the objective function provided in Trajectron++ [2] and WWTG [7], Where we aim to solve:

$$\begin{aligned} L_{reg} &= \max_{\phi, \theta, \psi} \sum_{r=1}^N \mathbb{E}_{z \sim q_{\phi}(z|V_{i_r}, V_{f_r}, M)} [\log p_{\psi}(V_{f_r}|V_{i_r}, z, M)] \\ &\quad - \beta D_{\text{KL}}(q_{\phi}(z|V_{i_r}, V_{f_r}, M) || p_{\theta}(z|V_{i_r}, M)) + \alpha I_q(V_{i_r}; z), \end{aligned} \quad (8)$$

where by [2], the I_q is the mutual information between V_i and z under the distribution $q_{\phi}(V_i, z)$. The computation of I_q is the same as [25]. α and β are hyperparameters.

Classification Loss: For the classification loss, we consider the methods in TPNNet [19], which a binary cross-entropy loss L_{class} is employed as:

$$L_{class}(x, y) = -w(y * \log(x) + (1 - y) * \log(1 - x)) \quad (9)$$

The total loss is written as follows:

$$L_{total} = L_{reg} + \frac{1}{N} \sum_n L_{class}(c_n, c_n^*) \quad (10)$$

where N is the number of trajectory proposals in the proposal set, w is learnable weight, and c_n^* is the corresponding predicted label, and c_n is the corresponding ground truth label.

During the training phase, the regression module minimizes the regression loss, and the classification module minimizes the classification loss.

V. EXPERIMENTS

A. Datasets

Our model is evaluated on four widely used public datasets: The ETH, UCY, Stanford Drone Dataset, and nuScenes. The ETH and UCY datasets focus on the pedestrian trajectory prediction, and contains complex social interactions. The ETH/UCY dataset has five subsets, each named ETH, HOTEL, UCY, ZARA-01, ZARA-02. There are two settings for the length of trajectories, $T_{obs} = T_{future} = 3.2s$ and $T_{obs} = 3.2s, T_{future} = 4.8s$. The data is captured at $2.5Hz$ ($\Delta t = 0.4s$), So the dataset will contains 8 frames for observations and 8/12 frames for prediction.

For the Stanford Drone Dataset, this is a trajectory dataset that is captured by drones from top-down view. So the scenes in the dataset are top-down-view. The scene are captured at a university campus with vehicles, cyclists, and crowds. The dataset contains a lot of heterogeneous data.

For the nuScenes dataset, this is a challenging large real-world driving forecasting dataset, where with more than 1000

scenes in the dataset are captured in Boston and Singapore. Each scene is 20 seconds long, and the dataset contains High-Definition semantic maps. All the scenes in the dataset contains a large amount of heterogeneous data, with complex social interactions among up to 23 semantic object classes. Also, the map provides data about the physical constraints in each scenes.

B. Evaluation Metrics

In our experiment, we will use four metrics that are commonly used in all trajectory prediction models, which include minADE_{20} , minFDE_{20} , ADE_{20} and FDE_{20} .

- minADE_{20} (mADE_{20}): the minimum mean ℓ_2 distance between the ground truth trajectory and predicted trajectory from the best of 20 samples by the LTN.
- minFDE_{20} (mFDE_{20}): the minimum ℓ_2 distance between the ground truth final position and the predicted final position at the final T_{future} from the best of 20 samples by the LTN.
- ADE: the mean ℓ_2 distance between the ground truth trajectory and predicted trajectory by the LTN.
- FDE: the ℓ_2 distance between the ground truth final position and the predicted final position at the final T_{future} by the LTN.

C. Baselines

We compare the performance of our model with the state-of-the-art models below:

- Vanilla LSTM: An LSTM network utilizing only the agent A_i 's history trajectory information.
- Social LSTM [4]: An LSTM network that utilizes not only the agent A_i 's history trajectory information, but also uses LSTM to model the agents' trajectory information around agent A_i .
- Social GAN [26] (S-GAN): This is a GAN with social interaction considered. Each agent is modeled by an LSTM-GAN combined network, where the LSTM encoder-decoder outputs are the generator of GAN, and the generated trajectories are then evaluated against the ground truth trajectories in the discriminator.
- Trajectron++ [2]: This is a CVAE-based trajectory prediction model, where the LSTM-GRU encoder-decoder structure is used, and the LSTM is used to model the agent's history and it's corresponding social interaction, and at the GRU decoder a full distribution of the predicted trajectory is produced.
- Social-BiGAT [9] (S-BiGAT): This is an LSTM-GAN with Graph Attention Network to encode agent's social interactions.
- TPNet [19]: This is a CNN based network that produces a trajectory proposal set by first predict the end point from the given map information and then predict the potential endpoint, then do regression based on map information, starting and the end point. At the end, a classification is performed to output proposals with high scores.

- Sophie [27]: This is a GAN-based trajectory prediction model that also leverages the social interactions and physical information. Similar to Social-GAN, the trajectory is produced by generator and the discriminator will evaluate these predictions against the ground truth trajectories.
- STGAT [28]: This is a spatial-temporal graph based trajectory prediction model. The spatial interaction is captured by graph attention mechanisms and the LSTM is used for temporal interactions.
- Social-Attention [29](S-ATTN): This is an attention based trajectory prediction model, where it captures the relative importance, which is attention, of each person when predicting for the future trajectories.
- Car-Net [30]: This is a prediction model that can account the dependencies between agent's behavior and their spatial environment, where the model can learn where to look in a large environment when predicting the trajectory of an agent.
- CSP [24]: This is a prediction model that builds on LSTM encoder-decoder framework, where LSTM is used to model the agent's social interaction and output multi-modal future distribution of the agent based on the social interaction.
- SpAGNN [31]: This is a probabilistic model that utilizes graph neural network to capture the interactions between the vehicles and output a distribution of the future trajectory of the vehicle the model selected to predict.

Implementation Details For our experiment, we adopted the ETH/UCY/nuScenes dataset preprocessing method provided in Trajectron++ [2]. For SDD, it is preprocessed according to the methods provided in Evolvegraph [32]. For our Mogrifier GRU, we chose the mogrifying step to be 6, and the α , β in the regression objective function is $\alpha=1$, and β is dynamically changed for the best performance, according to Trajectron++, which this methods provides most optimal result. We optimize the network using Adam [33] optimizer with learning rate 0.002. The γ we used is $3m$. We implemented LTN with PyTorch on a desktop with Ubuntu 20.04, equipped with one Intel I7-8700K CPU and two Nvidia GTX 1070Ti GPUs.

D. Evaluation of Trajectory Prediction

We compared our Mogrifier GRU version LTN with several baseline method on ETH, UCY, datasets in terms of two metrics ADE and FDE in Table I. We reported our results as LTN_{M20} , where it's the 20 trajectory proposal that has the highest scores. In Table I, we can see that in terms of ADE, our LTN_{M20} can outperform Trajectron++, which leads the rest of the model in the data we collected. In ETH and UNIV, our model achieved better ADE results. Also, in terms of FDE, our LTN_{M20} outperforms all other methods completely, showing that our model is very good at long term trajectory prediction. We successfully improved our methods by nearly 10% in terms of FDE, or long term prediction. However, we also notice that both ETH and UCY has reached a saturation, which means that more improvement is unlikely due to the data annotation errors or any off-errors during data collection, which leads us to analyze the LTN's performance in another dataset SDD.

TABLE I: Comparison between Mogrifler GRU version LTN (LTN_{M20}) and the baseline methods on ETH and UCY benchmark. We used $T_{future} = 4.8$ seconds setting. Each row represents a dataset and each column represents a method, with ADE and FDE separated into two tables for better clarity. The ADE and FDE are measured in Meter. All the $_{20}$ marks means that the prediction is chosen from the best of 20 samples.

Metric	Dataset	LSTM	S-LSTM	S-GAN ₂₀	Trajectron++ ₂₀	S-BiGAT ₂₀	TPNet ₂₀	SoPhie ₂₀	STGAT	LTN_{M20}
ADE	ETH	1.09	1.09	0.81	0.43	0.69	0.84	0.70	0.65	0.39
	HOTEL	0.86	0.79	0.81	0.12	0.69	0.24	0.76	0.49	0.16
	UNIV	0.61	0.67	0.72	0.22	0.4	0.42	0.54	0.55	0.20
	ZARA1	0.41	0.47	0.60	0.17	0.55	0.33	0.30	0.30	0.18
	ZARA2	0.52	0.56	0.34	0.12	0.30	0.26	0.38	0.36	0.15
	Average	0.70	0.72	0.42	0.20	0.36	0.42	0.54	0.48	0.22
Metric	Dataset	LSTM	S-LSTM	S-GAN ₂₀	Trajectron++ ₂₀	S-BiGAT ₂₀	TPNet ₂₀	SoPhie ₂₀	STGAT	LTN_{M20}
FDE	ETH	2.41	2.35	1.52	0.86	1.29	1.73	1.43	1.12	0.80
	HOTEL	1.91	1.76	1.61	0.19	1.01	0.46	1.67	0.66	0.18
	UNIV	1.31	1.40	1.26	0.43	1.32	0.94	1.24	1.10	0.41
	ZARA1	0.88	1.00	0.63	0.32	0.62	0.75	0.63	0.69	0.29
	ZARA2	1.11	1.17	0.84	0.25	0.75	0.60	0.78	0.60	0.23
	Average	1.52	1.54	1.18	0.41	1.00	0.90	1.15	1.08	0.38

TABLE II: Comparison between Mogrifler GRU version LTN (LTN_{M20}) and the baseline methods on SDD benchmark. Each row represents a dataset and each column represents a method, with miniADE and miniFDE separated into two tables for better clarity. All the miniADE and miniFDE are measured in Pixels and are the measurement for predictions at $T_{future} = 4.8$ seconds. All the $_{20}$ marks means that the prediction is chosen from the best of 20 samples.

Metric	Dataset	S-LSTM	S-GAN	S-ATTN	STGAT	Trajectron++	CAR-Net	LTN_{M20}
mADE ₂₀	SDD	31.4	27.0	33.3	18.8	19.3	25.72	15.2
Metric	Dataset	S-LSTM	S-GAN	S-ATTN	STGAT	Trajectron++	CAR-Net	LTN_{M20}
mFDE ₂₀	SDD	55.6	43.9	55.9	31.3	32.7	51.80	25.8

TABLE III: Comparison between Mogrifler GRU version LTN (LTN_{M20}) and the baseline methods on nuScenes benchmark. We predicted the trajectories 4 seconds after the T_{obs} . Each row represents methods, and each column represents the FDE of each methods at that discrete time period. All the FDE are measured in Meters.

Methods	1s	2s	3s	4s
S-LSTM	0.47	-	1.61	-
CSP	0.46	-	1.50	-
CAR-Net	0.38	-	1.35	-
SpAGGN	0.36	-	1.23	-
Trajectron++	0.07	0.45	1.14	2.20
LTN_{M20}	0.13	0.43	0.92	1.74

We compared our Mogrifler GRU version LTN with several baseline method on SDD dataset in terms of minADE₂₀ and minFDE₂₀ in table 2. In table 2, our methods outperformed all the baseline methods, with over 20 percent improvement in both minADE₂₀ and minFDE₂₀.

Not only our LTN demonstrated its long-term prediction advantages in terms of accuracy in ETH, UCY, and SDD, we also introduced a dataset with more heterogeneous data, nuScenes.

In Table III, we compared our LTN_{M20} model with several baseline methods with nuScenes dataset. We can see that Trajectron++ outperformed all the other baseline methods across all 4 seconds time span, but our LTN_{M20} outperforms Trajectron++, especially in long terms, where over 20 percent improvement is achieved in 3seconds and 4seconds. But we also noticed that in short term, especially in 1seconds, our method did not outperform Trajectron++, but we are not especially concerned with it because our model demonstrated persuasive

TABLE IV: Comparison between Mogrifler GRU version Trajectron++ ($Trajectron++_M$) and Regular GRU version Trajectron++ on nuScenes benchmark. We predicted the trajectories 4 seconds after the T_{obs} . All the FDE are measured in Meters.

Methods	1s	2s	3s	4s
Trajectron++	0.07	0.45	1.14	2.20
$Trajectron++_M$	0.08	0.43	1.02	2.03

long term prediction performance in terms of accuracy.

E. Evaluation of Mogrifler GRU

We now examine the performance of the Mogrifler GRU. As we mentioned before, because we believe both ETH/UCY datasets have reached a saturation, where there are no possible improvement space for us, we conduct the rest of the experiment on SDD and nuScenes datasets. As the table shows, we set up experiment with two versions of Trajectron++, since the authors kindly provided their code:

Version 1: Trajectron++ with regular GRU as decoder.

Version 2: Trajectron++ with Mogrifler GRU as decoder.

As Table IV shows, the Mogrifler GRU version Trajectron++ outperforms the regular GRU version Trajectron++ at 3 seconds and 4 seconds. At 3 seconds and 4 seconds the improvement reached around 10 percent just by switching the regular GRU with our Mogrifler GRU. Our Mogrifler GRU is intended to remind the future hidden states and input with the information before, so our experiment showed that with the reminder of the states of the agent in 1 seconds and 2 seconds provided to 3 seconds and 4 seconds time period, the prediction result in 3 seconds and 4 seconds improved.

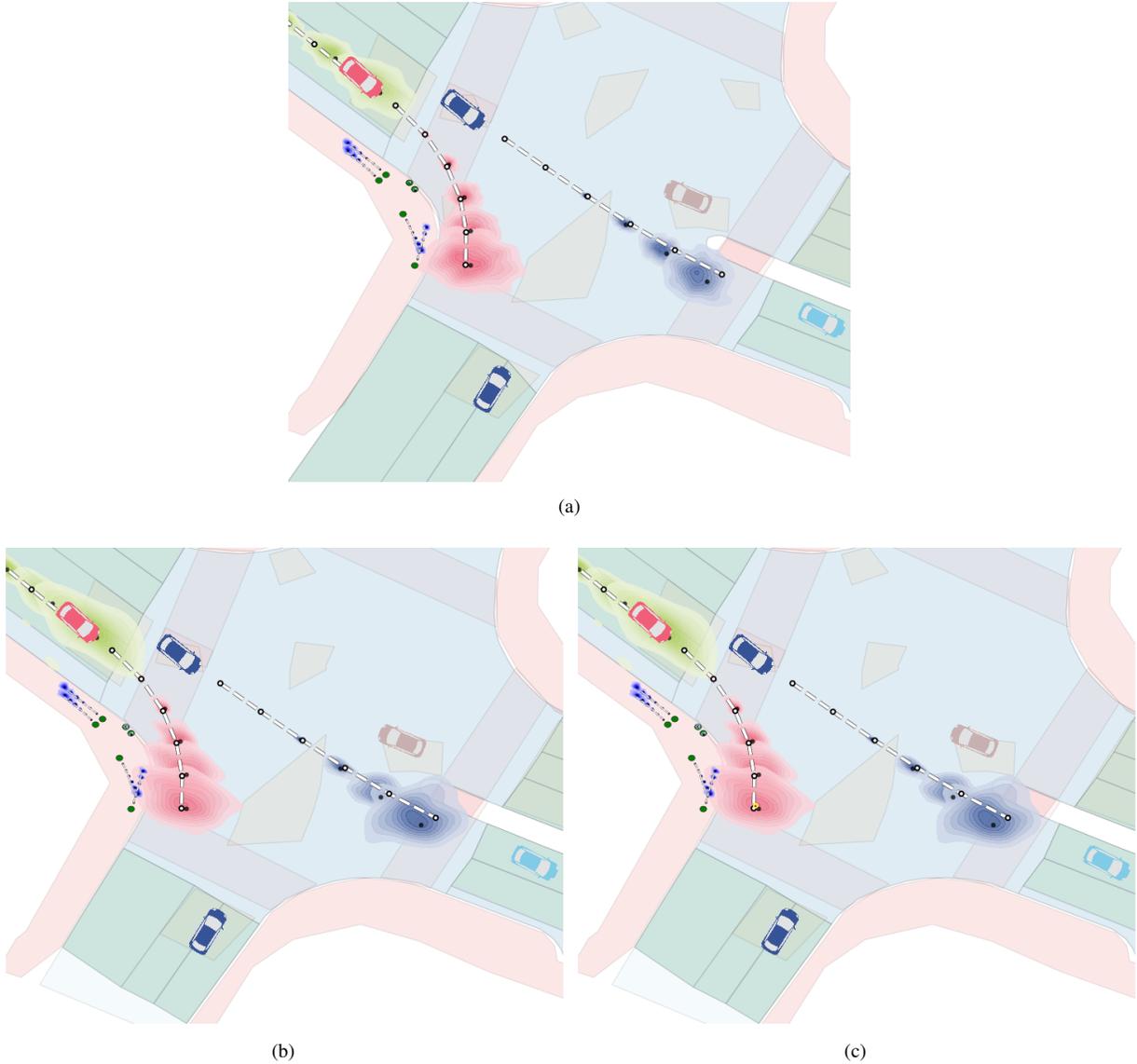


Fig. 2: The same scene with predictions from Trajectron++ [2] and LTN_{M20} . (a) The prediction from Trajectron++ where a most likely single output and the full distribution of future trajectory are provided. (b) The z_{mode} distribution prediction from LTN_{M20} and the most likely single output from the Trajectron++ are provided. (c) The prediction of $T_{future} = 3$ seconds from LTN_{M20} , denoted in a yellow circle, and the z_{mode} distribution are provided.

F. Qualitative Analysis

We compared our visualized experiment results with Trajectron++ [2] in the same scene provided in nuScenes [17] dataset. In Fig. 2(a) and Fig. 2(b), we can clearly see the advantage brought by the z_{mode} distribution. We can compare the shape of the distribution in both figures, where we can see that z_{mode} distribution of the red and blue car looks more uniform, meaning that it can generate samples with less bias. Therefore, the z_{mode} distribution in Fig. 2(b) can provide more samples near the ground truth. Also, the dense center of the z_{mode} distribution shifts towards the ground truth comparing with that of the full distribution. So, the proposed samples in z_{mode} will have larger probability to be the closest samples to the ground truth. Although the variance increases for the z_{mode} distribution, we will tackle it with our classification module,

where the samples with high variance will be filtered out by the threshold γ .

In Fig. 2(b), we can also see that by plainly taking the argmax of the z_{mode} distribution, the output is not the closest one to the ground truth. This is very clear in Fig. 2(b), where for the red car and its end point prediction, the z_{mode} indicates an area at the very center of the distribution that is closest to the ground truth trajectory, but the most likely output does not fall into that area. By our classification module, with good sampling from the z_{mode} distribution, filtering, and classification, we obtained the yellow hollow circle in Fig. 2(c) that falls into the area indicated by the z_{mode} distribution. The result is better compared with Trajectron++. For clarity, we only indicate the prediction of LTN_{M20} at $T_{future} = 3$ seconds, which is the end point distribution in red for the red

car.

VI. CONCLUSION

In this work, we present our model LTN, a two-stage trajectory prediction model for long-term trajectory prediction. The LTN incorporates heterogeneous data and combines regression and classification methods to improve the trajectory prediction performance for long-term prediction. Our LTN will first generate a distribution of the future trajectories, sample future trajectory proposals from it and perform classification on our proposal set. Along with the data of surrounding agents and the map information, our model could ensure that the final trajectory prediction is ideal and better. Our LTN achieves state-of-the-art performance in various popular datasets and shows significant improvement in terms of long-term trajectory prediction accuracy.

REFERENCES

- [1] J. Li, H. Ma, and M. Tomizuka, "Conditional generative neural system for probabilistic trajectory prediction," in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2019, pp. 6150–6156.
- [2] T. Salzmann, B. Ivanovic, P. Chakravarty, and M. Pavone, "Trajectron++: Dynamically-feasible trajectory forecasting with heterogeneous data," 2020.
- [3] W. Zhan, J. Li, Y. Hu, and M. Tomizuka, "Safe and feasible motion generation for autonomous driving via constrained policy net," in *IECON 2017-43rd Annual Conference of the IEEE Industrial Electronics Society*. IEEE, 2017, pp. 4588–4593.
- [4] 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 (CVPR)*, June 2016.
- [5] J. Li, H. Ma, and M. Tomizuka, "Interaction-aware multi-agent tracking and probabilistic behavior prediction via adversarial learning," in *2019 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2019.
- [6] H. Ma, J. Li, W. Zhan, and M. Tomizuka, "Wasserstein generative learning with kinematic constraints for probabilistic interactive driving behavior prediction," in *2019 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2019, pp. 2477–2483.
- [7] P. Felsen, P. Lucey, and S. Ganguly, "Where will they go? predicting fine-grained adversarial multi-agent motion using conditional variational autoencoders," in *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018.
- [8] W. Zhan, L. Sun, Y. Hu, J. Li, and M. Tomizuka, "Towards a fatality-aware benchmark of probabilistic reaction prediction in highly interactive driving scenarios," in *2018 21st International Conference on Intelligent Transportation Systems (ITSC)*. IEEE, 2018, pp. 3274–3280.
- [9] V. Kosaraju, A. Sadeghian, R. Martín-Martín, I. Reid, S. H. Rezatofighi, and S. Savarese, "Social-bigat: Multimodal trajectory forecasting using bicycle-gan and graph attention networks," 2019.
- [10] K. Sohn, H. Lee, and X. Yan, "Learning structured output representation using deep conditional generative models," in *Advances in Neural Information Processing Systems 28*, C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, Eds. Curran Associates, Inc., 2015, pp. 3483–3491.
- [11] J. Li, H. Ma, Z. Zhang, and M. Tomizuka, "Social-wagdat: Interaction-aware trajectory prediction via wasserstein graph double-attention network," *arXiv preprint arXiv:2002.06241*, 2020.
- [12] K. Cho, B. van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, "Learning phrase representations using rnn encoder-decoder for statistical machine translation," 2014.
- [13] G. Melis, T. Kočiský, and P. Blunsom, "Mogriifier lstm," 2020.
- [14] S. Pellegrini, A. Ess, K. Schindler, and L. van Gool, "You'll never walk alone: Modeling social behavior for multi-target tracking," in *2009 IEEE 12th International Conference on Computer Vision*, 2009, pp. 261–268.
- [15] L. Leal-Taixé, M. Fenzi, A. Kuznetsova, B. Rosenhahn, and S. Savarese, "Learning an image-based motion context for multiple people tracking," in *2014 IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 3542–3549.
- [16] A. Robicquet, A. Sadeghian, A. Alahi, and S. Savarese, "Learning social etiquette: Human trajectory understanding in crowded scenes," in *Computer Vision – ECCV 2016*, B. Leibe, J. Matas, N. Sebe, and M. Welling, Eds. Cham: Springer International Publishing, 2016, pp. 549–565.
- [17] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liong, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom, "nuscenes: A multimodal dataset for autonomous driving," *arXiv preprint arXiv:1903.11027*, 2019.
- [18] J. Li, H. Ma, W. Zhan, and M. Tomizuka, "Coordination and trajectory prediction for vehicle interactions via bayesian generative modeling," in *2019 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2019, pp. 2496–2503.
- [19] L. Fang, Q. Jiang, J. Shi, and B. Zhou, "Tpnet: Trajectory proposal network for motion prediction," 2020.
- [20] J. Li, H. Ma, W. Zhan, and M. Tomizuka, "Generic probabilistic interactive situation recognition and prediction: From virtual to real," in *2018 21st International Conference on Intelligent Transportation Systems (ITSC)*. IEEE, 2018, pp. 3218–3224.
- [21] J. Li, W. Zhan, Y. Hu, and M. Tomizuka, "Generic tracking and probabilistic prediction framework and its application in autonomous driving," *IEEE Transactions on Intelligent Transportation Systems*, vol. 21, no. 9, pp. 3634–3649, 2020.
- [22] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial networks," 2014.
- [23] B. Ivanovic, K. Leung, E. Schmerling, and M. Pavone, "Multimodal deep generative models for trajectory prediction: A conditional variational autoencoder approach," in *arXiv*, 2020.
- [24] N. Deo and M. M. Trivedi, "Multi-modal trajectory prediction of surrounding vehicles with maneuver based lstms," in *2018 IEEE Intelligent Vehicles Symposium (IV)*, 2018, pp. 1179–1184.
- [25] S. Zhao, J. Song, and S. Ermon, "Infovae: Balancing learning and inference in variational autoencoders," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, pp. 5885–5892, 07 2019.
- [26] A. Gupta, J. Johnson, L. Fei-Fei, S. Savarese, and A. Alahi, "Social gan: Socially acceptable trajectories with generative adversarial networks," 2018.
- [27] A. Sadeghian, V. Kosaraju, A. Sadeghian, N. Hirose, H. Rezatofighi, and S. Savarese, "Sophie: An attentive gan for predicting paths compliant to social and physical constraints," in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 1349–1358.
- [28] Y. Huang, H. Bi, Z. Li, T. Mao, and Z. Wang, "Stgat: Modeling spatial-temporal interactions for human trajectory prediction," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.
- [29] A. Vemula, K. Muelling, and J. Oh, "Social attention: Modeling attention in human crowds," 2018.
- [30] A. Sadeghian, F. Legros, M. Voisin, R. Vesel, A. Alahi, and S. Savarese, "Car-net: Clairvoyant attentive recurrent network," 2018.
- [31] S. Casas, C. Gulino, R. Liao, and R. Urtasun, "Spatially-aware graph neural networks for relational behavior forecasting from sensor data," 2019.
- [32] J. Li, F. Yang, M. Tomizuka, and C. Choi, "Evolvegraph: Multi-agent trajectory prediction with dynamic relational reasoning," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- [33] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," 2017.