

# Exploring Dynamic Context for Multi-path Trajectory Prediction

Hao Cheng<sup>1,\*</sup>, Wentong Liao<sup>2,\*</sup>, Xuejiao Tang<sup>2</sup>, Michael Ying Yang<sup>3</sup>, Monika Sester<sup>1</sup>, and Bodo Rosenhahn<sup>2</sup>

**Abstract**—To accurately predict future positions of different agents in traffic scenarios is crucial for safely deploying intelligent autonomous systems in the real-world environment. However, it remains a challenge due to the behavior of a target agent being affected by other agents dynamically and there being more than one socially possible paths the agent could take. In this paper, we propose a novel framework, named Dynamic Context Encoder Network (DCENet). In our framework, first, the spatial context between agents is explored by using self-attention architectures. Then, the two-stream encoders are trained to learn temporal context between steps by taking the respective observed trajectories and the extracted dynamic spatial context as input. The spatial-temporal context is encoded into a latent space using a Conditional Variational Auto-Encoder (CVAE) module. Finally, a set of future trajectories for each agent is predicted conditioned on the learned spatial-temporal context by sampling from the latent space, repeatedly. DCENet is evaluated on one of the most popular challenging benchmarks for trajectory forecasting *Trajnet* and reports a new state-of-the-art performance. It also demonstrates superior performance evaluated on the benchmark *inD* for mixed traffic at intersections. A series of ablation studies is conducted to validate the effectiveness of each proposed module. Our code is available at <https://github.com/wtliao/DCENet>.

## I. INTRODUCTION

Intelligent autonomous systems, such as robots and autonomous vehicles, have a high demand for the ability to accurately perceive, understand and predict the future behavior of humans for effective and safe deployments in our real-world environment. For example, an autonomous agent will adjust its moving path according to the possible locations of other agents to prevent obstructions or collisions. However, it is challenging to predict the future location of an agent because it is not deterministic: (1) an agent may change its mind during the movement, (2) other agents' behaviors will affect its next step (e.g., to avoid collisions), and (3) the influence from other agents is dynamic. Therefore, it is more beneficial to predict a set of potential trajectories adaptive to the dynamic interactions between agents than to predict a deterministic one. In this work, we seek to explore the dynamic context between agents in traffic scenarios to predict multiple possible trajectories for each agent in the short

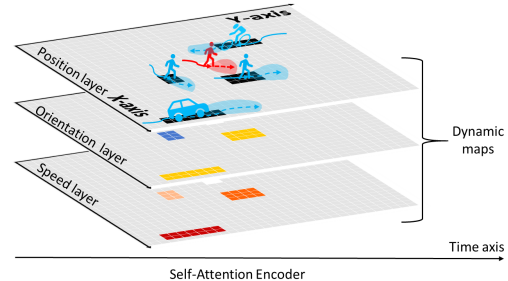


Fig. 1: Predicting multiple future trajectories (the most-likely one indicated by dash line over multiple ones indicated by shadow area) of a target agent (in red) conditioned on its observed movement (solid line) with the consideration of its interactions between neighboring agents (in blue) in mixed traffic. Interaction is learned through a sequence of dynamic maps at each step over the time axis and three layers are dedicated to capturing position, orientation and speed information (indicated by color-coded rectangles) using the self-attention structure.

future (12 steps) by observing their trajectories (8 steps), as showcased in Fig 1.

Specifically, the main contributions of this work are as follows: (1) It provides a novel framework to predict trajectories of heterogeneous agents (pedestrians, bicycles, vehicles, etc.) in various traffic situations, i.e., 20 different shared spaces and four intersections with mixed traffic. (2) Self-attention modules are integrated into our framework to explore the dynamic context among agents. (3) A set of possible trajectories for each agent is predicted conditioned on its observed trajectory and the learned dynamic context using a CVAE [1, 2] module. Extensive experiments are conducted on two of the most popular benchmarks *Trajnet* challenge [3] and the new large-scale benchmark *inD* [4] to validate the effectiveness of DCENet for trajectory forecasting. To judge the effectiveness of each proposed module, we conduct additional ablation studies. An overview of our framework is depicted in Fig. 2.

## II. RELATED WORK

**Trajectory Prediction.** Forecasting human trajectory has been researched for decades. In the early stages, many classic approaches are widely applied such as linear regression and Kalman filter [5], Gaussian processes [6] and Markov decision processing [7, 8]. These traditional methods heavily rely on the quality of manually designed features, which cannot work reliably in a real-world environment of complex

\*Equal contribution, name in alphabet order

<sup>1</sup>Institute of Cartography and Geoinformatics, Leibniz University Hannover, Germany, {cheng, sester}@ikg.uni-hannover.de

<sup>2</sup>Institute of Information Processing, Leibniz University Hannover, Germany, {lastname}@tnt.uni-hannover.de

<sup>3</sup>Scene Understanding Group, University of Twente, The Netherlands, michael.yang@utwente.nl

This work is supported by the German Research Foundation (DFG) through the Research Training Group SocialCars (GRK 1931) and Germany's Excellence Strategy within the Cluster of Excellence PhoenixD (EXC 2122).

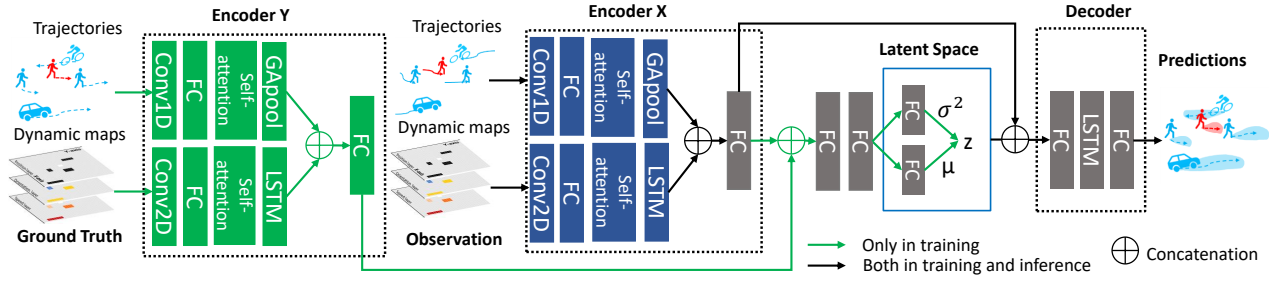


Fig. 2: The pipeline for the proposed method. The Encoder Y and Encoder X are identical in structure.

spatial-temporal dynamics and are poor at scaling up for dealing with a large amount of data. In recent years, many artificial intelligent (AI) technologies have been boosted by the cutting-edge deep learning technologies [9], including human trajectory prediction [10]–[16]. The deep learning models, especially Recurrent Neural Networks (RNNs) with Long Short-Term Memories (LSTMs), show great power in modeling complex social interactions between agents for collision avoidance and exploiting the time dependency for predicting futures [17]. The Social LSTM network [10] explores the interactions between pedestrians by connecting neighboring LSTMs in the social pooling layer and predicts trajectories for multiple pedestrians. Zhang *et al.* [13] propose the States Refinement LSTM (SR-LSTM) model that aligns all the agents together and refines the state of each agent through a message-passing framework. Chandra *et al.* [18] combine LSTM and Convolutional Neural Network (CNN) to model the interactions between heterogeneous road agents. However, many works have figured out the limited capability of LSTMs in modeling human-human interactions [19, 20]. Hence, the attention module [21] is incorporated in LSTMs to learn the spatial-temporal context of trajectories between pedestrians in [12, 22, 23]. Recently, the Transformer structure [24] has shown its power in context learning and sequential prediction [25, 26]. In this paper, we will adopt the self-attention module to encode the dynamic interactions between agents. The recent work [27] seeks to utilize the Transformer structure to predict trajectory instead of LSTMs. Our work is different from it essentially: (1) we use the generic self-attention module rather than the Deep Bidirectional Transformers (BERT) [25], which is a heavy stacked Transformer structure and is pre-trained on large-scale datasets, and (2) our framework is a generative model.

**Multi-path Trajectory Prediction.** Many approaches have been proposed to predict a socially compliant set of possible trajectories for an agent [11, 28]–[33]. Generative Adversarial Nets (GAN) [34] and CVAE [1, 2] are the most popular generative models used for this task. In [11] a trajectory sampler named Social GAN is proposed that considers the social effects of all agents. The generator is trained to predict a set of trajectories for each agent against a recurrent discriminator. In [12] social and physical attention mechanisms are implemented in the GAN sampler to predict paths for each agent. In [28], multiple plausible prediction samples are generated by a CVAE-based RNN encoder-

decoder conditioned on observations. Katyal *et al.* [33] propose to predict the intent of the target agent using a Bayesian approach as a condition of their CVAE-based LSTM encoder-decoder to help generate multiple paths. Meanwhile, they introduce an LSTM discriminator to train the framework in an adversarial way. Salzmann *et al.* [35] propose a CVAE-based model using spatial-temporal graphs to predict pedestrian and car trajectories. In [36], scene context and the interactions between individual and group agents are accounted as a condition in a CVAE-based framework to sample multiple trajectories. [37] applies a determinantal point process to increase the diversity sampling of a CVAE-based model for 2D and 3D motion prediction using synthetic data. Some other works treat the multi-path trajectory prediction problem as the estimation of a multimodal distribution. Cui *et al.* [32] propose to model the multimodality of vehicle movement prediction with Deep Convolutional Networks. In [30], first, the multimodal distributions are predicted with an evolving strategy by combining the Winner-Takes-ALL loss [38]. Then, the samples from the first stage fit a distribution for trajectory prediction. Cheng *et al.* [39] propose AMENet that only employs the self-attention mechanism [24] for learning agent-to-agent interaction. In comparison, DCENet adopts a two-stream architecture [40, 41] of attention modules, with respective streams dedicated to learning the spatial and temporal contexts explicitly.

### III. METHOD

#### A. Problem Formulation

Trajectory prediction is defined as to sequentially predict the future positions  $\hat{\mathbf{Y}}_i = \{\hat{\mathbf{y}}_i^{T+1}, \dots, \hat{\mathbf{y}}_i^{T'}\}$  of target agent  $i$  by observing its trajectory  $\mathbf{X}_i = \{\mathbf{x}_i^1, \dots, \mathbf{x}_i^T\}$ , where  $\mathbf{x}_i^t = (x_i^t, y_i^t)$  is the coordinates at the  $t$ -th step and  $1 \leq t \leq T$ . Similarly,  $\hat{\mathbf{y}}_i^{t'} = (x_i^{t'}, y_i^{t'})$  is the coordinates at the  $t'$ -th step and  $T < t' \leq T'$ .  $T$  is the length of observed trajectory and  $T'$  is the total length of being observed and predicted trajectory in discrete time steps.  $\hat{\mathbf{Y}}_i$  should be as close to the corresponding ground truth  $\mathbf{Y}_i$  as possible. The problem of multi-path trajectory prediction can be formulated as predicting a set of trajectories  $\hat{\mathbf{Y}}_i = \{\hat{\mathbf{Y}}_{i,1}, \dots, \hat{\mathbf{Y}}_{i,N}\}$  by observing  $\mathbf{X}_i$  for agent  $i$ , where  $N$  is the total number of predicted trajectories.

#### B. Dynamic Maps

To model the interactions among agents, we first create dynamic maps for each agent that consist of the orientation,

speed and position layers of its intermediate environment. These dynamic maps are different from the ones in [41] that are designed for modeling map rasterization and traffic lights. Centralized on the target agent, a map is defined as a rectangular area of size  $W \times H$  and divided into grid cells. First, referring to the target agent  $i$ , the neighboring agents  $N(i)$  are mapped into the closest grid cells  $s_{w \times h}^t$  according to their relative position as well as the cells reached by their anticipated relative offset (speed) in the  $x$  and  $y$  directions:

$$\begin{aligned} \text{cells}_w^t &= x_j^t - x_i^t + (\Delta x_j^t - \Delta x_i^t), \\ \text{cells}_h^t &= y_j^t - y_i^t + (\Delta y_j^t - \Delta y_i^t), \end{aligned} \quad (1)$$

where  $w \leq W$ ,  $h \leq H$ ,  $j \in N(i)$  and  $j \neq i$ . The *orientation layer*  $O$  stores the heading direction that is defined as the angle  $\vartheta_j$  in the Euclidean plane and calculated in the given radians by  $\vartheta_j = \arctan2(\Delta y_j^t, \Delta x_j^t)$ .  $(\Delta y_j^t, \Delta x_j^t)$  is the offset of the position from  $t$ -th step to the next one for neighboring agent  $j$ . The angle is shifted into degree  $[0, 360)$ . Similarly, the *speed layer*  $S$  stores the travel speed and the *position layer*  $P$  stores the position using a binary flag in the cells mapped above. Last, layer-wise, a Min-Max normalization scheme is applied for normalization, see Fig. 1. The map should cover a large vicinity area. Empirically we found  $32 \times 32 m^2$  a proper setting considering both the coverage and the computational cost. The cell size is set to  $1 \times 1 m^2$  as a balance to avoid the overlap of multiple agents in one cell based on the distribution of the experimental data, which is also supported by the preservation of personal space [42].

### C. Encoder Network

The spatial-temporal context from both the observation time and prediction time are encoded by Encoder X and Y, respectively. Both encoders have the same two-stream structure: both streams consist of stacked self-attention layers; as illustrated in Fig. 2 one stream is followed by a global average pooling (GApool), while the other one is followed by an LSTM module. The upper stream is trained to learn motion information from the observed trajectory, whose input is the locations vector of the observed trajectory of the target agent  $\mathbf{X}_i = \{\mathbf{x}_i^1, \dots, \mathbf{x}_i^T\} \in \mathbb{R}^{T \times 2}$ . The lower stream is trained to explore dynamic interactions among agents from the dynamic maps noted as  $DM = \{O, S, P\} \in \mathbb{R}^{T \times H \times W \times 3}$  (discussed in Sec. III-B). For simplicity, we take the upper stream for illustration. To get a sparse high dimensional representation,  $\mathbf{X}_i$  is first passed to a 1D convolution layer (Conv1D) and a fully connected (FC) layer. Each of them is followed by a ReLU non-linear activation. We denote this operation as  $\pi(\mathbf{X}_i)$ . A self-attention layer takes as input the Query ( $Q$ ), Key ( $K$ ) and Value ( $V$ ) and outputs a weighted sum of the value vectors. The weight assigned to each value is calculated as the dot-product of the query with the corresponding key:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V, \quad (2)$$

where  $\sqrt{d_k}$  is the scaling factor,  $d_k$  is the dimension of the vector  $K$  and  $T$  is the transpose operation. This operation is

also called *scaled dot-product attention* [24]. The  $Q$ ,  $K$  and  $V$  are obtained by three separated linear transformations:

$$Q = \pi(\mathbf{X})W_Q, \quad K = \pi(\mathbf{X})W_K, \quad V = \pi(\mathbf{X})W_V, \quad (3)$$

where  $W_Q, W_K, W_V \in \mathbb{R}^{d_\pi \times d_k}$  are the trainable parameters and  $d_\pi$  is the dimension of  $\pi(\mathbf{X})$ .

Because the self-attention module takes all inputs at the same time, position encodings are added to the  $Q$ ,  $K$  and  $V$  at the bottom of each self-attention layer to encode the temporal information. The sine and cosine functions of different frequencies (varying in time here) are the most widely used:

$$\mathbf{p}^t = \{p_{t,d}\}_{d=1}^D, \quad p_{t,d} = \begin{cases} \sin\left(\frac{t}{10000^{d/D}}\right), & \text{for } d \text{ even;} \\ \cos\left(\frac{t}{10000^{d/D}}\right), & \text{for } d \text{ odd,} \end{cases} \quad (4)$$

where  $D = d_k$  ensures position encodings to have the same dimension as the vectors of  $Q$ ,  $K$  and  $V$ .

To attend to different information from different representation subspaces jointly, the *multi-head attention* [24] strategy is applied as a conventional operation, where a head is an independent scaled dot-product attention module:

$$\begin{aligned} \text{MultiHead}(Q, K, V) &= \text{ConCat}(\text{head}_1, \dots, \text{head}_h)W_O, \\ \text{head}_i &= \text{Attention}(QW_{Qi}, KW_{Ki}, VW_{Vi}), \end{aligned} \quad (5)$$

where  $W_{Qi}, W_{Ki}, W_{Vi} \in \mathbb{R}^{D \times d_{ki}}$  are the linear transformation parameters same as in Eq. (3) and  $W_O$  are the linear transformation parameters for aggregating the extracted information from different heads. Note that  $d_{ki} = \frac{d_k}{h}$  and  $d_{ki}$  must be an aliquot part of  $d_k$ .  $h$  is the total number of the attention heads and we use two heads in the implementation.

Then the GApool is used to extract the temporal dependencies between steps by taking as input the output of the self-attention module and output an encoded representation.

The lower stream that exploits the dynamic interactions among agents works in the same way but the spatial dependencies among agents are encoded by the hidden states of an LSTM. Finally, the outputs of these two streams are connected and passed to a FC layer for fusion as the encoded information that includes dynamic spatial-temporal context.

### D. Multiple Trajectories Prediction

Our method is CVAE-based and predicts multiple trajectories by repeatedly sampling from a learned latent space conditioned on the encoded information. The CVAE is an extension of the VAE [43] by introducing a condition to control the output [2]. Given a set of samples  $(\mathbf{X}, \mathbf{Y}) = ((\mathbf{X}_1, \mathbf{Y}_1), \dots, (\mathbf{X}_m, \mathbf{Y}_m))$ , it jointly learns a recognition model  $q_\phi(\mathbf{z}|\mathbf{Y}, \mathbf{X})$  of a variational approximation of the true posterior  $p_\theta(\mathbf{z}|\mathbf{Y}, \mathbf{X})$  and a generation model  $p_\theta(\mathbf{Y}|\mathbf{X}, \mathbf{z})$  for predicting the output  $\mathbf{Y}$  conditioned on the input  $\mathbf{X}$ .  $\mathbf{z}$  are the stochastic latent variables,  $\phi$  and  $\theta$  are the respective recognition and generative parameters. The goal is to maximize the *Conditional Log-Likelihood*:  $\log p_\theta(\mathbf{Y}|\mathbf{X}) = \log \sum_{\mathbf{z}} p_\theta(\mathbf{Y}, \mathbf{z}|\mathbf{X}) = \log \left( \sum_{\mathbf{z}} q_\phi(\mathbf{z}|\mathbf{X}, \mathbf{Y}) \frac{p_\theta(\mathbf{Y}|\mathbf{X}, \mathbf{z}) p_\theta(\mathbf{z}|\mathbf{X})}{q_\phi(\mathbf{z}|\mathbf{X}, \mathbf{Y})} \right)$ . According to Jensen's inequality [44], the evidence lower bound

can be obtained:

$$\log p_\theta(\mathbf{Y}|\mathbf{X}) \geq -D_{KL}(q_\phi(\mathbf{z}|\mathbf{X}, \mathbf{Y})||p_\theta(\mathbf{z})) + \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{X}, \mathbf{Y})}[\log p_\theta(\mathbf{Y}|\mathbf{X}, \mathbf{z})], \quad (6)$$

where  $p_\theta(\mathbf{z})$  is made statistically independent from  $p_\theta(\mathbf{z}|\mathbf{X})$  [1, 2]. Here both the approximated posterior  $q_\phi(\mathbf{z}|\mathbf{X}, \mathbf{Y})$  and the prior  $p_\theta(\mathbf{z})$  are assumed to be Gaussian distribution for an analytical solution [43]. During training, the Kullback-Leibler divergence  $D_{KL}(\cdot)$  acts as a regularizer and pushes the approximated posterior to the prior distribution  $p_\theta(\mathbf{z})$ . The generation error  $\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{X}, \mathbf{Y})}(\cdot)$  measures the distance between the generated output and the ground truth. During inference, for a given observation  $\mathbf{X}_i$ , one latent variable  $\mathbf{z}_i$  is drawn from the prior distribution  $p_\theta(\mathbf{z})$ , and one of the possible output  $\hat{\mathbf{Y}}_i$  is generated from the distribution  $p_\theta(\mathbf{Y}_i|\mathbf{X}_i, \mathbf{z}_i)$ . The latent variables  $\mathbf{z}$  allow for the one-to-many mapping from the condition to the output via multiple sampling. In this work, we model a conditional distribution  $p_\theta(\mathbf{Y}_n|\mathbf{X})$ , where  $\mathbf{X}$  is the observed trajectory information and  $\mathbf{Y}_n$  is one of its possible future trajectories.

**Training:** As shown in Fig. 2, during the training, both the observed trajectory  $\mathbf{X}_i$  and its future trajectory  $\mathbf{Y}_i$  are encoded by Encoder X and Y (see Sec. III-C), respectively. Then, their encodings are concatenated and passed through two FC layers (each is followed by a ReLU activation) for fusion. Then, two side-by-side FC layers are used to estimate the mean  $\mu_{z_i}$  and the standard deviation  $\sigma_{z_i}$  of the latent variables  $\mathbf{z}_i$ . A trajectory  $\hat{\mathbf{Y}}_i$  is reconstructed by an LSTM decoder step by step by taking  $\mathbf{z}_i$  and the encodings of observation as input. Because the random sampling process of  $\mathbf{z}_i$  can not be back propagated during training, the standard reparameterization trick [43] is adopted to make it differentiable. To minimize the error between the predicted trajectory  $\hat{\mathbf{Y}}_i$  and the ground truth  $\mathbf{Y}_i$ , the reconstruction loss is defined as the L2 loss (Euclidean distance). Thus, the whole network is trained by minimizing the loss function using the stochastic gradient descent method:

$$L = \|\hat{\mathbf{Y}} - \mathbf{Y}\|^2 + D_{KL}(q_\phi(\mathbf{z}|\mathbf{X}, \mathbf{Y})||\mathcal{N}(0, I)). \quad (7)$$

**Test:** In the test phase, the ground truth of future trajectory is no more available and its pathway is removed (color coded in green in Fig. 2). A latent variable  $\mathbf{z}$  is sampled from the prior distribution  $\mathcal{N}(0, I)$  and concatenated with the observation encodings that serve as the condition for the following trained decoder, so that the decoder can predict a trajectory. To predict multiple trajectories, this process (sampling and decoding) is repeated multiple times.

### E. Trajectory Ranking

We propose a ranking strategy to select the *most-likely* predicted trajectory out of the multiple predictions in order to adjust the Trajnet challenge setting. We apply bivariate Gaussian distribution to rank the predicted trajectories  $(\hat{\mathbf{Y}}_{i,1}, \dots, \hat{\mathbf{Y}}_{i,N})$  for each agent. At step  $t'$ , all the predicted positions for agent  $i$  are stored in  $|\hat{\mathbf{X}}_i, \hat{\mathbf{Y}}_i|^{t'}$ . We follow [45]

to fit the positions into the probability density function:

$$f(\hat{x}_i, \hat{y}_i)^{t'} = \frac{1}{2\pi\sigma_{\hat{x}_i}\sigma_{\hat{y}_i}\sqrt{1-\rho^2}} \exp \frac{-Z}{2(1-\rho^2)},$$

$$Z = \frac{(\hat{x}_i - \mu_{\hat{x}_i})^2}{\sigma_{\hat{x}_i}^2} + \frac{(\hat{y}_i - \mu_{\hat{y}_i})^2}{\sigma_{\hat{y}_i}^2} - \frac{2\rho(\hat{x}_i - \mu_{\hat{x}_i})(\hat{y}_i - \mu_{\hat{y}_i})}{\sigma_{\hat{x}_i}\sigma_{\hat{y}_i}}. \quad (8)$$

where  $\mu$  denotes the mean and  $\sigma$  the standard deviation, and  $\rho$  is the correlation between  $\hat{X}_i$  and  $\hat{Y}_i$ . A predicted trajectory is scored as the sum of the relative likelihood of all its steps:  $S(\hat{\mathbf{Y}}_{i,n}) = \sum_{t'=T+1}^{T'} f(\hat{x}_i, \hat{y}_i)^{t'}$ . All predicted trajectories are ranked by this score and the one with the highest score stands out for the single-path prediction.

## IV. EXPERIMENTS

To evaluate the performance of our proposed method, we compare DCENet with the most influential and recent nine state-of-the-art models from the Trajnet [3] challenge leader-board for a fair comparison: (1) *Linear (off)*: a simple temporal linear regressor; (2) *Social Force* [46]: the very high impact rule-based model that implements social force to avoid collisions; (3) *S-LSTM* [10]: the highly cited LSTM-based model that introduces social pooling layer for modeling interactions; (4) *S-GAN* [11]: a GAN-based trajectory predictor; (5) *MX-LSTM* [47]: an LSTM trajectory predictor that utilizes the head direction of agent; (6) *SR-LSTM* [13]: an LSTM-based model that refines the hidden states by message passing; (7) *RED* [19]: an RNN encoder-decoder model predicts trajectory only using observations; (8) *Ind-TF* [27]: a Transformer-based trajectory predictor; (9) *AMENet* [39]: the most recent state-of-the-art on the Trajnet leader-board. We further design a series of ablation studies to analyze the impact of each proposed module, *i.e.*, dynamic maps, transformer and LSTM encoder/decoder: (1) *Baseline*: an LSTM encoder-decoder only using the observed trajectory as input; (2) *DCENet w/o DMs*: the stream of encoding dynamic maps is removed from our final model; (3) *Trans. En&De*: the LSTM encoder-decoder is substituted by the Transformer encoder/decoder [24] in our framework.

### A. Datasets

**Trajnet** [3] is one of the most popular forecasting benchmarks. In Trajnet, 8 consecutive ground-truth locations (3.2 seconds) of each trajectory are for observation and the following 12 steps (4.8 seconds) are required to forecast. Trajnet is a superset of diverse popular benchmark datasets: ETH [48], UCY [49], Stanford Drone Dataset [50], BIWI Hotel [48], and MOT PETS [51]. There is a total of 11448 trajectories from these four subsets covering 38 scenes for training. The test data is from the diverse partitions of them (besides MOT PETS) of the other 20 scenes without ground truth. The Trajnet challenge provides a specific server for online evaluation. It is worth noting that many existing works are evaluated on a subset of Trajnet using their own train/test splits. For the sake of fairness, we only compare DCENet to the works which have shown their performance on the Trajnet challenge leader-board.

TABLE I: Results of different methods on the Trajnet challenge [3]. Models are categorized into deterministic (determ.) and stochastic (stoch.) depending on whether they incorporate a generative module.

Model	Category	Avg. [m]↓	FDE [m]↓	ADE [m]↓
S-LSTM [10]	determ.	1.3865	3.098	0.675
S-GAN [11]	stoch.	1.3340	2.107	0.561
MX-LSTM [47]	determ.	0.8865	1.374	0.399
Linear (off)	determ.	0.8185	1.266	0.371
Social Force [46]	determ.	0.8185	1.266	0.371
SR-LSTM [13]	determ.	0.8155	1.261	0.370
RED [19]	determ.	0.7800	1.201	0.359
Ind-TF [27]	determ.	0.7765	1.197	0.356
AMENet [39]	stoch.	0.7695	1.183	0.356
Baseline	stoch.	0.8045	1.239	0.370
DCENet w/o DMs	stoch.	0.7760	1.195	0.357
Trans. En&De	stoch.	0.7780	1.196	0.360
DCENet	stoch.	<b>0.7660</b>	<b>1.179</b>	<b>0.353</b>

**inD** was acquired by Bock *et al.* [4] using drones at four busy intersections in Germany in 2019. The traffic is dominated by vehicles and they interact with pedestrians heavily. The speed difference and confrontation makes the trajectory prediction challenging. The data was processed to obtain the same format as Trajnet: 8 steps for observation and the following 12 steps for prediction.

### B. Evaluation Metrics

We adopt the most popular evaluation metrics: the mean average displacement error (ADE) and the final displacement error (FDE) to measure the trajectory prediction performance. ADE measures the aligned Euclidean distance from the prediction to its corresponding ground truth trajectory averaged over all steps. The mean value across all the trajectories is reported. FDE measures the Euclidean distance between the last position from the prediction to the corresponding ground truth position. In addition, the most-likely prediction is decided by the ranking method as described in Sec III-E. Compared with the ground truth (only if it is available), *@top10* is the one out of ten predicted trajectories that has the smallest ADE and FDE.

The implementation details of training and testing our methods can be found in our code repository.

### C. Results

The experimental results from different methods including our ablative models reported on the Trajnet leader-board are listed in Table I. Without ground truth trajectories, the single-path trajectory prediction was selected by the ranking mechanism. We can see that DCENet reported new state-of-the-art performance and the ablative models also had comparable performances compared to the previous works.

First, by comparing to the Baseline, both DCENet w/o DMs and Ind-TF had much better results, and DCENet w/o DMs was slightly better in the average score and FDE but a little inferior in ADE than Ind-TF. Considering both models only use observed trajectories as input, it indicates that our method (self-attention + LSTM encoder/decoder) explored

TABLE II: Quantitative results of our model and the comparative models on the inD benchmark measured by ADE/FDE.

Model	S-LSTM	S-GAN	AMENet	DCENet
inD	@ <i>top 10</i>			
Intersection-(A)	2.04/4.61	2.84/4.91	0.95/1.94	<b>0.72/1.50</b>
Intersection-(B)	1.21/2.99	1.47/3.04	0.59/1.29	<b>0.50/1.07</b>
Intersection-(C)	1.66/3.89	2.05/4.04	0.74/1.64	<b>0.66/1.40</b>
Intersection-(D)	2.04/4.80	2.52/5.15	0.28/0.60	<b>0.20/0.45</b>
Avg.	1.74/4.07	2.22/4.29	0.64/1.37	<b>0.52/1.23</b>
inD	<i>Most-likely</i>			
Intersection-(A)	2.29/5.33	3.02/5.30	1.07/2.22	<b>0.96/2.12</b>
Intersection-(B)	1.28/3.19	1.55/3.23	0.65/1.46	<b>0.64/1.41</b>
Intersection-(C)	1.78/4.24	2.22/4.45	<b>0.83/1.87</b>	0.86/1.93
Intersection-(D)	2.17/5.11	2.71/5.64	0.37/0.80	<b>0.28/0.62</b>
Avg.	1.88/4.47	2.38/4.66	0.73/1.59	<b>0.69/1.52</b>

a better spatial-temporal context than Transformer. Furthermore, Ind-TF utilizes BERT, a heavily stacked Transformer structure and must be pre-trained on an external large-scale dataset, while DCENet does not require it. The results of DCENet w/o DMs indicates that its superior performance is not because we used more information (dynamic maps).

Second, by the comparison between the Baseline and S-LSTM, our Baseline model was significantly better. The difference between them is that our Baseline is CVAE-based and generates multiple trajectories. It indicates that the future motion of humans is of high uncertainty, and predicting a set of possible trajectories is better than only predicting a single one. It also demonstrates the effectiveness of the trajectory ranking methods (Sec. III-E), which was used to select the most-likely trajectory from the multiple predictions. Our Baseline outperformed S-GAN significantly, which is a generative model for multiple trajectories prediction.

Third, interestingly, Trans. En&De that adopts the Transformer encoder and decoder in our framework did not achieve improved performance compared to DCENet. This phenomenon indicates that our self-attention + LSTM encoder/decoder structure explored better dynamic context between agents than Transformer encoder/decoder in terms of trajectory prediction. The superior performance of DCENet w/o DMs against Ind-TF has also confirmed that.

Lastly, DCENet outperformed DCENet w/o DMs. It indicates that the dynamic maps helped model the interactions between agents and were useful for trajectory prediction.

**Discussion** According to the comparison above, the results indicate: (1) DCENet is effective for predicting accurate trajectories for heterogeneous agents in various real-world traffic scenes, even without modeling interactions explicitly (the Baseline model). (2) The ranking method correctly estimates the multiple predictions and recommends a reliable candidate for the single-path trajectory prediction task. (3) Compared to the Baseline model, DCENet learns interaction via the dynamic maps with the self-attention structure effectively and shows improved performance. (4) Both LSTM and Transformer networks are capable of learning complex sequential patterns but their combination further enhances the performance in terms of trajectory prediction.



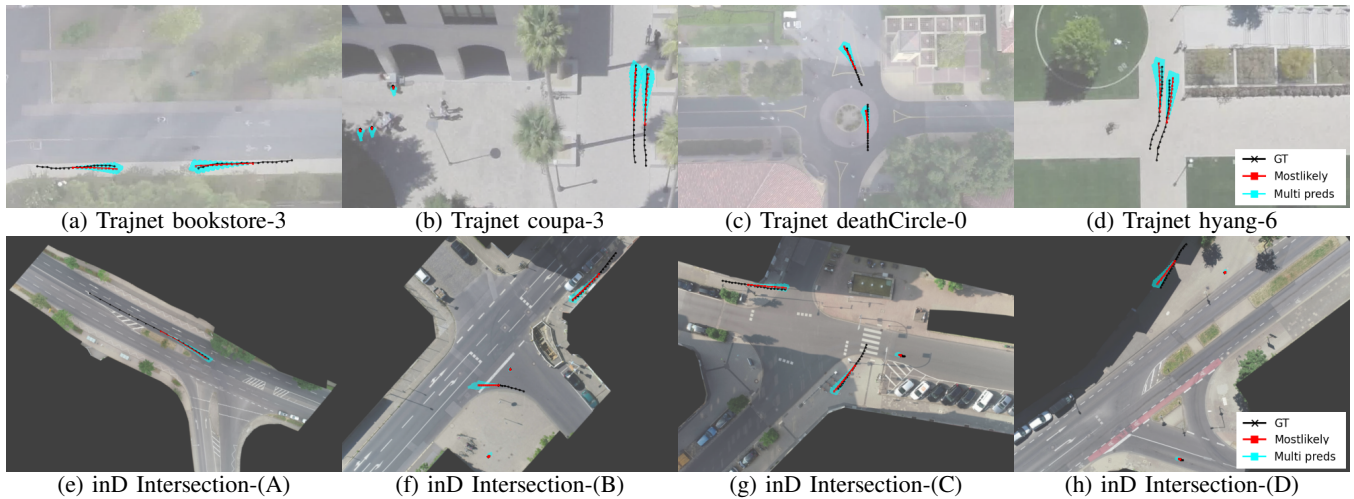


Fig. 3: Multi-path trajectory predictions in shared spaces in Trajnet (1st row) and at different intersections in inD (2nd row).

Furthermore, we have tested DCENet on inD [4] to justify its performance and generalization ability. We compare our model with the three most relevant models: S-LSTM for comparing with its occupancy grid mapping for agent-to-agent interaction, S-GAN for its generative module, and AMENet for its CVAE module and LSTM sequential modeling. To guarantee a fair comparison, all the models were trained and tested using the same data. S-LSTM predicts the distributions of the positions [10]. During inference, multiple positions were generated by sampling. Table II lists the performance measured by ADE/FDE. Our model achieved the best performance for the  $@top10$  prediction across all the intersections and reduced the errors by a big margin. Our model also outperformed the other models for the most-likely prediction at three out of four intersections. It only slightly fell behind the AMENet model on the intersection-(C). We anticipate that the most-likely prediction fell behind the  $@top10$  prediction. However, the ranking method was still effective in recommending a reliable candidate in comparison to the other models. The results indicate: (1) Our model is able to generalize on different datasets and maintain superior performance. (2) Predicting multiple paths is more beneficial than predicting a single one for an agent. On the one hand, multiple predictions increase the chances to narrow down the errors. On the other hand, a single prediction may lead to a wrong conclusion especially if the initial steps predicted are deviating from the ground truth and the errors will accumulate significantly with time. The multiple predictions form into an area indicating the potential intent of an agent and the area size reflects the uncertainty of an agent's intent.

The qualitative results are shown in Fig. 3. The first row showcases the scenarios in the Trajnet dataset. Note that the qualitative analysis on Trajnet was carried out on the validation set (an independent subset of the training set) for comparing with the ground truth. Our model accurately predicted two pedestrians walking towards each other at bookstore-3. The shadow areas indicate multiple possible trajectories. It also correctly predicted the static pedestrians

in coupa-3, as well as the pedestrians walking in parallel. In deathCircle-0, our model predicted different possible turning angles for the cyclist in the roundabout. In hyang-6, two pedestrians walking closely to each other were predicted correctly. The second row showcases the scenarios in the inD dataset. Our model predicted a fast driving vehicle with a slightly different predicted speed at the Intersection-(A). It predicted that a left-turning vehicle may turn at the intersection-(B) with varying tuning angle and speed. The model also correctly predicted the interaction at the zebra crossing at the intersection-(C), where the vehicle stops to yield the way to the pedestrian. Similar predictions can be seen for the walking and static pedestrians, as well as the vehicle waiting at the entrance of the intersection-(D). Overall, we can also see that the recommended single path is very close to the corresponding ground truth for each agent.

## V. CONCLUSION

In this paper, we proposed a novel framework DCENet for multi-path trajectory prediction for heterogeneous agents in various real-world traffic scenarios. We decompose the learning of dynamic spatial-temporal context into exploiting the dynamic spatial context between agents using self-attention and the LSTM encoder and learning temporal context between steps with the following self-attention and global average pooling. The spatial-temporal context is encoded into a latent space using a CVAE module. Finally, a set of future trajectories for each agent is predicted conditioned on the spatial-temporal context using the trained CVAE module. DCENet was evaluated on the Trajnet challenge benchmark and achieved the new state-of-the-art performance on the leader-board. Its superior performance on the inD benchmark further validated its efficacy and generalization ability. The ablation studies justified the impact of each module in DCENet. In the future, we are interested in extending the method for learning the impact from environment/static context, *e.g.*, space layout and scene deployment, to further enhance the performance of trajectory prediction.

## REFERENCES

- [1] D. P. Kingma, S. Mohamed, D. J. Rezende, and M. Welling, “Semi-supervised learning with deep generative models,” in *NeurIPS*, 2014, pp. 3581–3589.
- [2] K. Sohn, H. Lee, and X. Yan, “Learning structured output representation using deep conditional generative models,” in *NeurIPS*, 2015, pp. 3483–3491.
- [3] A. Sadeghian, V. Kosaraju, A. Gupta, S. Savarese, and A. Alahi, “Trajnet: Towards a benchmark for human trajectory prediction,” *arXiv preprint*, 2018.
- [4] J. Bock, R. Krajewski, T. Moers, S. Runde, L. Vater, and L. Eckstein, “The ind dataset: A drone dataset of naturalistic road user trajectories at german intersections,” *arXiv preprint arXiv:1911.07602*, 2019.
- [5] A. C. Harvey, *Forecasting, structural time series models and the Kalman filter*. Cambridge university press, 1990.
- [6] M. K. C. Tay and C. Laugier, “Modelling smooth paths using gaussian processes,” in *Field and Service Robotics*, 2008, pp. 381–390.
- [7] D. Makris and T. Ellis, “Spatial and probabilistic modelling of pedestrian behaviour,” in *BMVC*, 2002, pp. 1–10.
- [8] K. M. Kitani, B. D. Ziebart, J. A. Bagnell, and M. Hebert, “Activity forecasting,” in *ECCV*, 2012, pp. 201–214.
- [9] Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning,” *Nature*, vol. 521, no. 7553, p. 436, 2015.
- [10] A. Alahi, K. Goel, V. Ramanathan, A. Robicquet, L. Fei-Fei, and S. Savarese, “Social lstm: Human trajectory prediction crowded spaces,” in *CVPR*, 2016, pp. 961–971.
- [11] A. Gupta, L. Johnson, Justand Fei-Fei, S. Savarese, and A. Alahi, “Social gan: Socially acceptable trajectories with generative adversarial networks,” in *CVPR*, 2018, pp. 2255–2264.
- [12] A. Sadeghian, V. Kosaraju, A. Sadeghian, N. Hirose, and S. Savarese, “Sophie: An attentive gan for predicting paths compliant to social and physical constraints,” in *CVPR*, 2019, pp. 1349–1358.
- [13] P. Zhang, W. Ouyang, P. Zhang, J. Xue, and N. Zheng, “Sr-lstm: State refinement for lstm towards pedestrian trajectory prediction,” in *CVPR*, 2019, pp. 12 085–12 094.
- [14] N. Mohajerin and M. Rohani, “Multi-step prediction of occupancy grid maps with recurrent neural networks,” in *CVPR*, 2019, pp. 10 600–10 608.
- [15] C. Tang and R. R. Salakhutdinov, “Multiple futures prediction,” in *NeurIPS*, 2019, pp. 15 398–15 408.
- [16] R. Chandra, U. Bhattacharya, A. Bera, and D. Manocha, “Traffic: Trajectory prediction dense and heterogeneous traffic using weighted interactions,” in *CVPR*, 2019, pp. 8483–8492.
- [17] P. Kothari, S. Kreiss, and A. Alahi, “Human trajectory forecasting crowds: A deep learning perspective,” *arXiv preprint arXiv:2007.03639*, 2020.
- [18] R. Chandra, T. Guan, S. Panuganti, T. Mittal, U. Bhattacharya, A. Bera, and D. Manocha, “Forecasting trajectory and behavior of road-agents using spectral clustering in graph-lstms,” *IEEE Robotics and Automation Letters*, vol. 5, no. 3, pp. 4882–4890, 2020.
- [19] S. Becker, R. Hug, W. Hübner, and M. Arens, “An evaluation of trajectory prediction approaches and notes on the trajnet benchmark,” *arXiv preprint arXiv:1805.07663*, 2018.
- [20] S. Becker, R. Hug, W. Hubner, and M. Arens, “Red: A simple but effective baseline predictor for the trajnet benchmark,” in *ECCV*, 2018, pp. 138–153.
- [21] D. Bahdanau, K. Cho, and Y. Bengio, “Neural machine translation by jointly learning to align and translate,” in *ICLR*, 2015.
- [22] A. Al-Molegi, M. Jabreel, and A. Martinez-Balleste, “Move, attend and predict: An attention-based neural model for people’s movement prediction,” *Pattern Recognition Letters*, vol. 112, pp. 34–40, 2018.
- [23] A. Vemula, K. Muelling, and J. Oh, “Social attention: Modeling attention human crowds,” in *ICRA*, 2018, pp. 1–7.
- [24] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” in *NeurIPS*, 2017, pp. 5998–6008.
- [25] J. Devlin, M. Chang, K. Lee, and K. Toutanova, “BERT: pre-training of deep bidirectional transformers for language understanding,” in *Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2019, pp. 4171–4186.
- [26] S. He, W. Liao, H. R. Tavakoli, M. Yang, B. Rosenhahn, and N. Pugeault, “Image captioning through image transformer,” in *ACCV*, 2020.
- [27] F. Giuliari, I. Hasan, M. Cristani, and F. Galasso, “Transformer networks for trajectory forecasting,” in *ICPR*, 2020.
- [28] N. Lee, W. Choi, P. Vernaza, C. B. Choy, P. H. Torr, and M. Chandraker, “Desire: Distant future prediction dynamic scenes with interacting agents,” in *CVPR*, 2017, pp. 336–345.
- [29] J. Amirian, J.-B. Hayet, and J. Pettré, “Social ways: Learning multimodal distributions of pedestrian trajectories with gans,” in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*. IEEE, 2019, pp. 2964–2972.
- [30] O. Makansi, E. Ilg, O. Cicek, and T. Brox, “Overcoming limitations of mixture density networks: A sampling and fitting framework for multimodal future prediction,” in *CVPR*, 2019, pp. 7144–7153.
- [31] A. Poibrenski, M. Klusch, I. Vozniak, and C. Müller, “M2p3: multimodal multi-pedestrian path prediction by self-driving cars with egocentric vision,” in *Annual ACM Symposium on Applied Computing*, 2020, pp. 190–197.
- [32] H. Cui, V. Radosavljevic, F.-C. Chou, T.-H. Lin, T. Nguyen, T.-K. Huang, J. Schneider, and N. Djuric, “Multimodal trajectory predictions for autonomous driving using deep convolutional networks,” in *ICRA*, 2019, pp. 2090–2096.
- [33] K. D. Katyal, G. D. Hager, and C.-M. Huang, “Intent-aware pedestrian prediction for adaptive crowd navigation,” in *ICRA*, 2020, pp. 3277–3283.
- [34] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative adversarial nets,” in *NeurIPS*, 2014, pp. 2672–2680.
- [35] T. Salzmann, B. Ivanovic, P. Chakravarty, and M. Pavone, “Trajectron++: Dynamically-feasible trajectory forecasting with heterogeneous data,” in *ECCV*, vol. 12363. Springer, 2020, pp. 683–700.
- [36] H. Cheng, W. Liao, M. Y. Yang, M. Sester, and B. Rosenhahn, “Mcnnet: Multi-context encoder network for homogeneous agent trajectory prediction mixed traffic,” in *ITSC*, 2020.
- [37] Y. Yuan and K. M. Kitani, “Diverse trajectory forecasting with determinantal point processes,” in *ICLR*, 2020.
- [38] A. Guzman-Rivera, D. Batra, and P. Kohli, “Multiple choice learning: Learning to produce multiple structured outputs,” in *NeurIPS*, 2012, pp. 1799–1807.
- [39] H. Cheng, W. Liao, M. Y. Yang, B. Rosenhahn, and M. Sester, “Amenet: Attentive maps encoder network for trajectory prediction,” *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 172, pp. 253–266, 2021.
- [40] K. Simonyan and A. Zisserman, “Two-stream convolutional networks for action recognition in videos,” in *NeurIPS*, 2014, pp. 568–576.
- [41] S. Casas, W. Luo, and R. Urtasun, “Intentnet: Learning to predict intention from raw sensor data,” in *Conference on Robot Learning*. PMLR, 2018, pp. 947–956.
- [42] C. L. Gérin-Lajoie, Martand Richards and B. J. McFadyen, “The negotiation of stationary and moving obstructions during walking: anticipatory locomotor adaptations and preservation of personal space,” *Motor control*, vol. 9, no. 3, pp. 242–269, 2005.
- [43] D. P. Kingma and M. Welling, “Auto-encoding variational bayes,” in *ICLR*, 2014.
- [44] J. L. W. V. Jensen *et al.*, “Sur les fonctions convexes et les inégalités entre les valeurs moyennes,” *Acta mathematica*, vol. 30, pp. 175–193, 1906.
- [45] A. Graves, “Generating sequences with recurrent neural networks,” *arXiv preprint arXiv:1308.0850*, 2013.
- [46] D. Helbing and P. Molnar, “Social force model for pedestrian dynamics,” *Physical review E*, vol. 51, no. 5, p. 4282, 1995.
- [47] I. Hasan, F. Setti, T. Tsesmelis, A. Del Bue, F. Galasso, and M. Cristani, “Mx-lstm: mixing tracklets and vislets to jointly forecast trajectories and head poses,” in *CVPR*, 2018, pp. 6067–6076.
- [48] S. Pellgrini, A. Ess, K. Schindler, and L. Van Gool, “You’ll never walk alone: Modeling social behavior for multi-target tracking,” in *ICCV*, 2009, pp. 261–268.
- [49] A. Lerner, Y. Chrysanthou, and D. Lischinski, “Crowds by example,” in *Computer Graphics Forum*, vol. 26, no. 3. Wiley Online Library, 2007, pp. 655–664.
- [50] A. Robicquet, A. Sadeghian, A. Alahi, and S. Savarese, “Learning social etiquette: Human trajectory understanding crowded scenes,” in *ECCV*, 2016, pp. 549–565.
- [51] J. Ferryman and A. Shahrokni, “Pets2009: Dataset and challenge,” in *International workshop on performance evaluation of tracking and surveillance*, 2009, pp. 1–6.