# Pedestrian Trajectory Prediction via Spatial Interaction Transformer Network

Tong Su, Yu Meng and Yan Xu

*Abstract*— As a core technology of the autonomous driving system, pedestrian trajectory prediction can significantly enhance the function of active vehicle safety and reduce road traffic injuries. In traffic scenes, when encountering with oncoming people, pedestrians may make sudden turns or stop immediately, which often leads to complicated trajectories. To predict such unpredictable trajectories, we can gain insights into the interaction between pedestrians. In this paper, we present a novel generative method named Spatial Interaction Transformer (SIT), which learns the spatio-temporal correlation of pedestrian trajectories through attention mechanisms. Furthermore, we introduce the conditional variational autoencoder (CVAE) [1] framework to model the future latent motion states of pedestrians. In particular, the experiments based on large-scale traffic dataset nuScenes [2] show that SIT has an outstanding performance than state-of-the-art (SOTA) methods. Experimental evaluation on the challenging ETH [3] and UCY [4] datasets confirms the robustness of our proposed model.

## I. INTRODUCTION

Vulnerable Roads Users(VRUs), due to their high maneuverability, may change their motion in a while. For the protection of pedestrians [5] [6], active vehicle safety systems (AVSSs) leverage the environmental perception and decision making technologies to minimize the effect of traffic accidents [7]. As a core component of AVSSs, pedestrian trajectory prediction module is responsible for providing warnings when pedestrians are close to the driving vehicles. By analyzing the movement patterns of other traffic agents [8] and predicting their future positions, vehicles equipped with prediction module are able to make appropriate navigation decisions (e.g. avoid impending collision) [9].

In chaotic traffic scenes, reliable trajectory prediction is challenging due to pedestrians react differently according to the change of surrounding environment. For example, pedestrians plan their future routes by sensing each other's posture or subtle changes in motion [10]. Therefore, the study of spatial interaction is essential for predicting future trajectories in the scenes with high pedestrian density.

To model multi-pedestrian trajectories in the scene, we build a dynamic graph to capture the complex spatial interaction. Different from previous works [11], transformer-based network is introduced to model the spatio-temporal dependencies. We believe that the powerful attention mechanisms are suitable for sequence modeling. Besides, we expect to aggregate the future latent motion states of pedestrians by their future trajectories. For example, in Fig. 1, pedestrian *P2* walks straight during the observation phase, but in the

The authors are with the School of Mechanical Engineering, University of Science and Technology Beijing, China. (e-mail: g20198575@xs.ustb.edu.cn; myu@ustb.edu.cn; b20160225@xs.ustb.edu.cn)
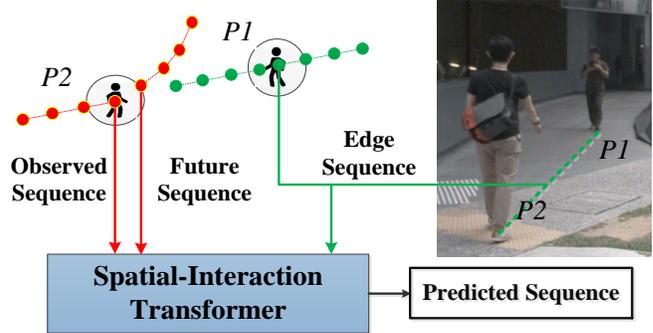
Fig. 1. Scene of interest, in where we output future positions of pedestrians with our proposed transformer-based methods that build a social graph to model the interaction between pedestrians. Furthermore, we consider the future motion states by incorporating their future sequence into our model.

forecast phase, *P2* make a sudden turn to avoid pedestrian *P1*. Only relying on observed trajectories is not enough to predict its future trajectories accurately, thus, we introduce the conditional variational autoencoder (CVAE) [1] framework that conditions future sequence to our prediction. With this framework, our model is able to make full use of labels to capture the future latent motion state.

In this work, we are interested in exploring the capability of transformer-based network in modeling social interaction. Specifically, our network tends to provide feasible approaches for pedestrian trajectory prediction in heavy traffic environments. The main contributions are as follows:

1): We present a novel deep generative model named spatial interaction transformer (SIT) that utilizes the attention mechanisms to dynamically model the spatial locations of the pedestrians and predict their future trajectories.

2): To handle the variance of pedestrian trajectories during observation and prediction, we use a generative framework that follows CVAE to incorporate pedestrian future trajectories into social interaction.

3): Extensive experiments are performed on two public datasets. Further statistical analyses show the effectiveness and robustness of our proposed data processing methods and edge modules.

The rest of the paper is organized as follows. In Section II, we introduce the related work about pedestrian trajectory prediction and describe the transformer network for trajectory prediction. Section III presents the detailed structure of our proposed SIT. The experiments and ablation study are performed in Section IV. Finally, Section V concludes our paper and provides the plans for future work.

## II. RELATED WORK

There are a large number of published studies that describe various pedestrian trajectory prediction methods. This section aims to focus on the literature relevant to our research. For this purpose, we introduce two aspects of related work, that are (a) pedestrian trajectory prediction and (b) transformer network for trajectory forecasting.

### A. Pedestrian Trajectory Prediction

**Traditional approaches**: Pedestrian trajectory prediction has attracted much attention in recent years [12]. Forecasting methods are mainly divided into two categories: kinetic-based forecasting methods and data-driven methods [13]. In [14], Schneider et al. used the extended Kalman filter and interactive multi-model to predict the diverse pedestrian motion states (stationary, interactive, bending, starting) and analyzed their differences. Since Long Short-Term Memory (LSTM) is more effective in modeling long-term sequence than kinetic-based methods, it has become the most frequently used model in trajectory prediction field. In [15], Li et al. compared the effects of different learning-based methods like Gaussian Process (GP), LSTM, GP-LSTM, Character-based LSTM, Sequence-to-Sequence (Seq2Seq) and attention-based Seq2Seq, and showed that the encoder-decoder structure such as Seq2Seq has an outstanding performance for modeling both linear or non-linear patterns.

**Spatial interaction**: The moving routes of pedestrians in traffic scenes are simply affected by surrounding agents. Some researchers have realized that by fusing the latent motion clues of surrounding pedestrians into trajectory prediction, it is possible to capture the dynamic changes of predicted trajectories. So far, social interaction has been extensively investigated by many works [16]. Various pedestrian trajectory analysis techniques have been proposed, ranging from deterministic linear regression [11] to generative model [17]. Social LSTM [11] encodes the historical trajectories of all pedestrians equally in the same scene and pools the obtained feature vectors by social pooling layer, which implicitly models the interaction of different pedestrians by learning spatial correlation. GRIP regards traffic objects as nodes and builds dynamic graphs to learn the movement patterns of pedestrians [18]. To predict the positions of all agents in the same scene at the same time, a spatio-temporal graph convolutional network is introduced to process the interaction of traffic subjects.

### B. Transformer Network for Trajectory Prediction

Although LSTM [15] has been applied in a variety of situations, it is difficult for LSTM to enhance its computational speed and performance due to its sequential structure. Transformer [19] network, as the SOTA models for most natural language processing tasks, can rely on powerful attention mechanisms to avoid these shortcomings. By feeding past positions into the network at the same timestep, transformer network has superior parallel computing capabilities and can learn concerning information from any historical position. As a result, transformer network has great potential to achieve remarkable performance in the field of pedestrian trajectory prediction.

Recently, Fran et al. [20] applied the transformer network to trajectory prediction. After inferencing in the velocity space, they used the predicted velocity to get the coordinate position of the pedestrians. However, they only considered the case of modeling a single pedestrian. In order to cope with crowd trajectory prediction, we propose a new transformer-based model to simulate the interaction between pedestrians.

## III. METHODOLOGY

In this section, we regard the trajectory prediction problem as a sequence regression task. Based on the historical observation, we aim to predict the trajectories of pedestrians in the future timesteps. Our proposed probabilistic generative model SIT is shown in Fig. 2.

### A. Input and Output of Model

**Data processing**: In order to learn the distribution of trajectory samples more effectively, processing the trajectory data is an essential part. The most commonly adopted trajectory processing method is to use the mean and standard deviation of all pedestrian positions to normalize the data [20]. However, in the different scenes, there is a large difference between the distributions of pedestrian trajectories. This method can not make predictions well in multiple scenarios [21]. Different from it, we use each pedestrian's last position of observed trajectories as the mean and the attention radius as the standard deviation to normalize each sample data separately, which makes the input data distribution more compact and improves the accuracy of prediction in various scenes.

**Input and output sequence**: Given the previous $H$ timesteps, we aim to predict trajectory horizon of $P$ timesteps with the time interval of $\triangle t$. At time $t$, the observed sequence is described as $X_{obs} = \{x^t\}_{t-H+1}^{t}$, future sequence is $Y = \{y^t\}_{t+1}^{t+P}$ and predicted sequence is $\hat{Y} = \{\hat{y}^t\}_{t+1}^{t+P}$. To make full use of abundant pedestrian motion state, we input 6-dimensional vectors including the normalized positions $(x_{pos}^t, y_{pos}^t)$, velocities $(x_{vel}^t, y_{vel}^t)$, and accelerations $(x_{acc}^t, y_{acc}^t)$. The predicted trajectories can generate more variability by adding velocity and acceleration feature to input data. Then, our model directly outputs the position $y^t$.

$$x^t = [x_{pos}^t, y_{pos}^t, x_{vel}^t, y_{vel}^t, x_{acc}^t, y_{acc}^t] \quad (1)$$
$$y^t = [x_{pos}^t, y_{pos}^t] \quad (2)$$
$$\hat{y}^t = [\hat{x}_{pos}^t, \hat{y}_{pos}^t] \quad (3)$$

**Positional encoding**: Unlike LSTM inputting the trajectory data step by step, SIT feeds into the input data simultaneously. Therefore, we manually add the timestamp information by using sine and cosine functions. Here d is
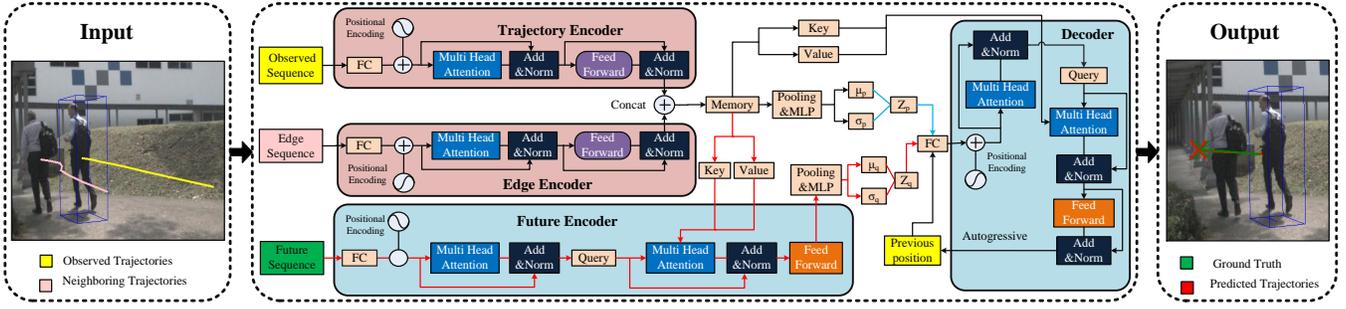
Fig. 2. The details of our proposed SIT. Based on CVAE, our model encodes observed sequence and their neighboring sequence to extract the temporal and spatial dependencies of pedestrians respectively. Red, blue, black arrows represent training only, inferencing only, and both of them, respectively.

each dimension of embedding.

$$P_{obs} = \begin{cases} sin(\frac{t}{10000^{d/D}}) & for \quad d \quad even \\ cos(\frac{t}{10000^{d/D}}) & for \quad d \quad odd \end{cases} \quad (4)$$

**Edge sequence**: In order to predict the positions of pedestrians accurately, we first build a graph $G = (V, E)$ to dynamically simulate the interaction between pedestrians and their neighbors. Each pedestrian is represented as a node $\nu \in V$ and edge $e = (\nu_i, \nu_j) \in E$ exist when $\nu_i$ influences $\nu_j$. In this work, we pay attention to the neighbors within the attention radius like other works [11] and euclidean distance indicate the edge $e$. For instance, the edge $e$ is taken into account when $\|p_i - p_j\| \leq$ attention radius where $p_i$, $p_j \in R^2$ are the 2D spatial position of $\nu_i, \nu_j$. In particular, the position $p_i$ is used to normalize the historical positions of neighbors $p_j$. This special technique can make all nodes in the scene have same latent states, which is beneficial to model the spatial interaction.

For every single pedestrian $\nu_i$, we merge its past states (normalized positions, velocities, and accelerations) of neighbors $\nu_j$ when edge $e$ is present, which can be achieved by element-wise sum. In this way, we convert the variable length neighbor states to a fixed edge sequence $X_{edge}$ which has the same shape as the observed sequence $X_{obs}$.

*B. Spatial Interaction Transformer Network*

The spatial interaction transformer network is mainly composed of attention mechanisms that assign unequal importance to neighboring pedestrians. Furthermore, our model follows the conditional variational autoencoder (CVAE) by modeling the future pedestrian trajectories as distributions based on their own and neighbors's past trajectories. We aim to learn the probabilistic model $p(Y|X_{obs}, X_{edge})$ by introducing latent variable $Z$. Here, $Z$ represents the latent state of pedestrian trajectories. We can describe the future trajectory distributions as the following equation:

$$p(Y|X_{obs}, X_{edge}) = \int p(Y|X_{obs}, X_{edge}, Z) \\ p(Z|X_{obs}, X_{edge})dZ \quad (5)$$

where $p(Z|X_{obs}, X_{edge})$ is the gaussian prior distribution which is inferred by past observed sequence $X_{obs}$ and

edge sequence $X_{edge}$. $p(Y|X_{obs}, X_{edge}, Z)$ is the conditional likelihood distribution that is impossible to calculate directly. So, to tackle this problem, we use the Kullback-Leibler divergence loss(KL loss) [1] as one item of our loss function:

$$L_{kl} = KL(q(Z|Y, X_{obs}, X_{edge}) \parallel p(Z|X_{obs}, X_{edge})) \quad (6)$$

where $q(Z|Y, X_{obs}, X_{edge})$ is the approximate posterior distribution. The latent state $Z$ can be jointly inferred by the future sequence $Y$, observed sequence $X_{obs}$ and edge sequence $X_{edge}$. This design allows $Z$ to consider not only its own future trajectories but also its neighboring trajectories, which enables our model to generate more interactive trajectories. $KL$ quantifies the difference between two probability distributions. By minimizing the $KL$, we can approximate the prior and posterior distributions. During training, we can infer $Z$ through $q(Z|Y, X_{obs}, X_{edge})$ and during testing $Z$ can be inferred by $p(Z|X_{obs}, X_{edge})$. After above formulation, we now introduce the details of our model.

**Trajectory and Edge Encoder**: We first add the edge sequence and observed sequence into the timestamp information through positional encoding and get the embedding vector. Both of them are fed into the trajectory encoder and edge encoder. Specifically, we encode the observed embeddings and edge embeddings with Multi-head attention and feed-forward network. After encoding, the concatenation operation is used to get the memory vector $C = \{c^t\}_{t-H+1}^t$ where $C$ summaries the past trajectories and the influence of all neighboring pedestrians. Then, a mean pooling layer is performed across all historical timesteps to get the past trajectory feature $c_n = mean(c_{t-H+1}, \ldots, c_t)$. We use a multi layer perception (MLP) to map $c_n$ to the gaussian prior distribution $p(Z|X_{obs}, X_{edge})$ and get the gaussian parameters $(\mu_p, \sigma_p)$. According to the Gumbel-Softmax reparameterization [1], we can sample $Z_p$ from the latent states:

$$Z_p = \mu_p\epsilon + \sigma_p, \quad \epsilon \sim N(0, 1) \quad (7)$$

**Future Encoder**: Given the pedestrian future trajectory sequence $Y$, we can obtain the timestamped sequence by positional encoding. After encoding by Multi-head attention, this sequence is feed into another Multi-head attention served

as queries. At the same time, the past trajectory memory $C$ is encoded as keys and values. The keys represent the weights for different timesteps and the values represent the latent state of different timesteps. This design allows our model to condition $X_{obs}$ through $C$, which is beneficial to approximate the posterior distribution effectively. Similar to prior distribution, a mean pooling layer is performed across future timesteps to extract the future feature and we use MLP to map future feature to the approximate posterior distribution $q(Z|Y, X_{obs}, X_{edge})$. Finally, we get the gaussian parameters $(\mu_q, \sigma_q)$ and sample $Z_q$ from those parameters:

$$Z_q = \mu_q \epsilon + \sigma_q, \quad \epsilon \sim N(0,1) \tag{8}$$

**Future Decoder**: It is worth noting that our future decoder is autoregressive, which means it outputs trajectory one step at a time. The input sequence of decoder can be described as $\{f^t\}_{t+1}^{t+P} = \{\hat{y}^t \oplus Z\}_{t+1}^{t+P}$ Here, $\hat{y}^{t+1}$ is initialized from the last position feature of observed sequence $X_{obs}$ and $Z$ is the sample $Z_p$ (training) or $Z_q$ (testing).

In the decoder stage, we add timestamp information into $f^t$ through positional encoding and feed them into the Multi-head attention. After obtaining the query vector, we input it into another Multi-head attention along with past trajectory memory served as keys and values. Then feed forward network is applied to output next timestep trajectory. Our future decoder allows the model to inference future trajectories while considering the effect of current neighbors. To approximate the conditional likelihood distribution $p(Y|X_{obs}, X_{edge}, Z)$ according to $q(Z|Y, X_{obs}, X_{edge})$, we minimize the mean squred error between predicted trajectories and future trajectory labels. Our loss function is written as:

$$L = \min \| Y - \hat{Y} \| + L_{kl} \tag{9}$$

### C. Training Details

We use batch size 100 on the training set and testing set. The latent states $|Z| = 32$. A 3-layer encoder and a 3-layer decoder are applied to our network. For data augmentation, we rotate the observed trajectories by an angle that varies from $0°$ to $360°$ with an interval of $15°$.

During training, we maintain the same weight initialization for all layers of SIT. The same as [20], we use 8 heads for Multi-head attention and D = 256. Following previous work [11], attention radius is set to 10 meters. We also use Adam as our optimizer with a decaying learning rate. Then the SIT is trained for 100 epochs.

### IV. EXPERIMENTS

In this section, we evaluate our proposed method on two public datasets. Besides, we perform the ablation study to quantitatively describe the effects of different components. It is worth noting that the experiments include varying prediction horizons.

### A. Datasets

We mainly use two public datasets to verify the performance of our proposed method. Both of them are in the world coordinate system. The first is the large-scale traffic dataset nuScenes for autonomous driving, which is annotated at 2 frames per second ($\triangle t = 0.5$). We extract scenes that contain pedestrians and get a training set of 632 scenes and a test set of 133 scenes. Each scene is 20s long. The total training sample contains 75767 sequences, and the test sample contains 12876 sequences. For the nuScenes dataset, we predict the trajectories of future 3s ($H$=6) based on the observed trajectories of past 4s ($P$=8). The second is the widely used benchmark datasets ETH (ETH and HOTEL) and UCY (UNIV, ZARA1, and ZARA2) in the field of pedestrian trajectory prediction. The datasets contain 5 different scenarios, each of which contains complex pedestrian interaction behaviors and is annotated at 2.5 Hz ($\triangle t = 0.4$). In order to maintain a fair comparison, we predict the trajectory horizon of 8 timesteps (3.2s) based on the observed length of 12 timesteps (4.8s) same as most papers [17].

### B. Evaluation Metrics

The same as prior papers [22], [23], we use the following metrics to evaluate our methods:

MAD (Mean Average Displacement, equivalently Average Displacement Error ADE): The average value of the Euclidean distance error between the predicted trajectories and the ground truth in the future time horizon T.

FAD (Final Average Displacement, equivalently Final Displacement Error FDE): The average value of the Euclidean distance error between the predicted trajectories and the ground truth at the last time step.

### C. NuScenes Dataset

**Experiment results**: To evaluate the effectiveness of our model, We compared the deterministic version of our proposed SIT with other baseline methods:

1) Constant Velocity(CV): This model assumes that pedestrians move at a constant velocity to linearly reason about the future trajectories.
2) LSTM: An LSTM encoder-decoder predicts future locations by using the past trajectories of pedestrians.
3) LSTM+Attention: Based on the LSTM encoder-decoder, an attention layer is added to the encoder. The attention mechanisms can automatically search for parts of the observed trajectories which are closely related to the the predicted sequence.
4) Transformer [20]: We keep the original model completely but use the nuScenes dataset for pedestrian trajectory prediction task.

Following the above methods, we separately conduct the experiments on the nuScenes dataset to compare the performance with prediction horizon varying from 1s to 3s. As shown in table I and Fig. 3, the prediction error MAD and FAD increase with the growth of time horizon. In addition, our model has lowest prediction errors over baseline methods

(outperforming existing approach [20] by 53% on MAD in the case of predicting 3s).

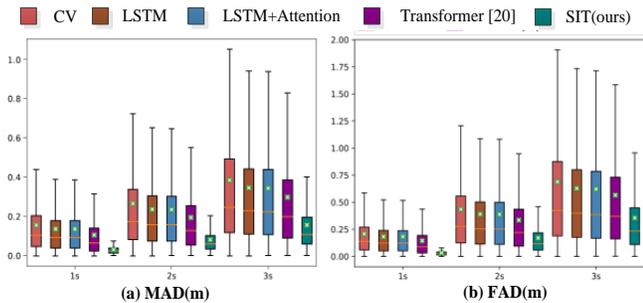| Methods | MAD(meters) | | | FAD(meters) | | |
|---|---|---|---|---|---|---|
| | 1s | 2s | 3s | 1s | 2s | 3s |
| Const.Velocity | 0.16 | 0.27 | 0.39 | 0.21 | 0.44 | 0.69 |
| LSTM | 0.14 | 0.24 | 0.35 | 0.18 | 0.39 | 0.63 |
| LSTM+Attention | 0.14 | 0.24 | 0.34 | 0.18 | 0.39 | 0.62 |
| Transformer [20] | 0.11 | 0.20 | 0.30 | 0.15 | 0.34 | 0.57 |
| **SIT(ours)** | **0.03** | **0.08** | **0.16** | **0.03** | **0.17** | **0.36** |



Fig. 3. The quantitative MAD and FAD results of all methods on the nuScenes test set when predicting 1-3s. We sample all test trajectories and use the boxplots to describe the distributions of their mean errors. "x" markers indicate the MAD or FAD value.

Then, we project the predicted trajectories to the image plane to visualize deterministic predicted examples of above mentioned prediction methods. As shown in Fig. 4, our proposed SIT performs particularly well in situations where a pedestrian may begin to turn suddenly. Our CVAE method plays a major role in modeling the nonlinear dynamics.

**Ablation study**: A comprehensive ablation study is performed in table II. We first make a qualitative analysis on the basis of the original transformer (the first row of table II) [20]. To fairly compare the impact of different components, we adopt the data processing method as described in III-A. The performance of original transformer data processing method can be seen in the fourth row of I. There is a slight improvement when adding to the edge encoding, which is crucial for modeling the spatial interaction. As can be seen in the third row, the CVAE framework yields a drastic reduction whether in MAD or FAD.

TABLE II

ABLATION STUDY

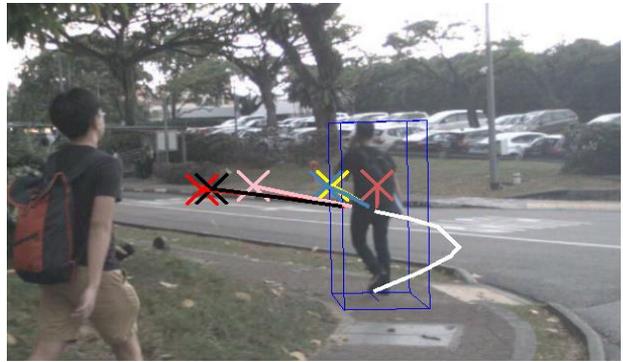| Components | | MAD | | | FAD | | |
|---|---|---|---|---|---|---|---|
| E[a] | C[b] | 1s | 2s | 3s | 1s | 2s | 3s |
| - | - | 0.07 | 0.13 | 0.22 | 0.09 | 0.24 | 0.46 |
| ✓ | - | 0.06 | 0.11 | 0.19 | 0.07 | 0.20 | 0.39 |
| ✓ | ✓ | **0.03** | **0.08** | **0.16** | **0.03** | **0.17** | **0.36** |

a. Edge encoding; b. CVAE



Fig. 4. Predicted trajectory examples for different baseline methods. The blue box represents the target pedestrian. The white line is the past positions and the red line is the ground truth of future positions. The other colors represent the respective trajectories predicted by different methods. We can see that SIT is closer to the ground truth.

*D. ETH and UCY Datasets*

To verify the robustness of our proposed method, we perform corresponding verifications on the benchmark datasets ETH and UCY in the prediction field. The leave-one-out evaluation strategy is generally adopted by most works. Specifically, we use 4 datasets of ETH and UCY to train the model and the rest to test.

As presented in table III, the deterministic model output one single trajectory and except ETH dataset our model has achieved outstanding performance among all SOTA methods. The performance of stochastic model is summaried in table IV. Here we sample 20 times and report the best sample. We can observe that our SIT significantly outperforms the baselines whether in MAD or FAD. One interesting finding is that our model significantly outperforms on UCY and ZARA, where crowd density is relatively higher. This can be explained by that our model is suitable to model the human-huamn interaction.

TABLE III

DETERMINISTIC MODEL

| Datasets | S-LSTM [11] | | TF[c] [20] | | STAR-D [24] | | **SIT(ours)** | |
|---|---|---|---|---|---|---|---|---|
| | M[a] | F[b] | M | F | M | F | **M** | **F** |
| ETH | 1.09 | 2.35 | 1.03 | 2.10 | **0.56** | **1.11** | 0.59 | 1.28 |
| HOTEL | 0.79 | 0.76 | 0.36 | 0.71 | 0.26 | 0.50 | **0.22** | **0.45** |
| UCY | 0.67 | 1.40 | 0.53 | 1.32 | 0.52 | 1.15 | **0.40** | **0.98** |
| ZARA1 | 0.47 | 1.00 | 0.44 | 1.00 | 0.41 | 0.90 | **0.30** | **0.75** |
| ZARA2 | 0.56 | 1.17 | 0.34 | 0.76 | 0.31 | 0.71 | **0.23** | **0.59** |
| avg[d] | 0.72 | 1.54 | 0.54 | 1.17 | 0.41 | 0.87 | **0.35** | **0.81** |

a. MAD; b. FAD; c. Transformer
d.The average value of 5 datasets

## V. DISCUSSION AND FUTURE WORK

This research aims to provide a novel method for pedestrian trajectory prediction task. Specifically, we propose deep

| Datasets | S-GAN [17] | | TF [20] | | STAR [24] | | SIT(ours) | |
|---|---|---|---|---|---|---|---|---|
| | M | F | M | F | M | F | **M** | **F** |
| ETH | 0.81 | 1.52 | 0.61 | 1.12 | **0.36** | **0.65** | 0.38 | 0.88 |
| HOTEL | 0.72 | 1.61 | 0.18 | 0.30 | 0.17 | 0.36 | **0.11** | **0.21** |
| UCY | 0.60 | 1.26 | 0.35 | 1.65 | 0.31 | 0.62 | **0.20** | **0.46** |
| ZARA1 | 0.34 | 0.69 | 0.22 | 0.38 | 0.26 | 0.55 | **0.16** | **0.37** |
| ZARA2 | 0.42 | 0.84 | 0.17 | 0.32 | 0.22 | 0.46 | **0.12** | **0.27** |
| avg | 0.58 | 1.18 | 0.31 | 0.55 | 0.26 | 0.53 | **0.19** | **0.44** |

generative model SIT which is based on attention mechanisms. In contrast to LSTM's sequential structure, our model learns the spatio-temporal correlation at a deeper level thanks to the self-attention network. We also merge the ground truth of future sequence into our trajectory prediction model. With CVAE, our model can make full use of future sequence. The effectiveness of our proposed modules is proved by our ablation study. On the other hand, we explore the ability of the transformer-based model to encode the spatial interaction between pedestrians. The results suggest that our method can achieve a significant improvement than SOTA methods [20].

Only adopting the pedestrian trajectories to make predictions, existing methods might fail in some complex scenes. In the development of the paper, we focus on fusing extra information like vehicles [25] and scenes [26] through transformer network. Moreover, future work will also pay attention to the variable length history trajectories because missing values [24] often exist in real world. We expect that the transformer network can solve this problem without padding or interpolation.

## ACKNOWLEDGMENT

## REFERENCES

[1] K. Sohn, H. Lee, and X. Yan, "Learning structured output representation using deep conditional generative models," *Advances in neural information processing systems*, vol. 28, pp. 3483–3491, 2015.

[2] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, and O. Beijbom, "nuscenes: A multimodal dataset for autonomous driving," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.

[3] 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.

[4] 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.

[5] D. Ferguson, M. Darms, C. Urmson, and S. Kolski, "Detection, prediction, and avoidance of dynamic obstacles in urban environments," in *2008 IEEE Intelligent Vehicles Symposium*. IEEE, 2008, pp. 1149–1154.

[6] C. Sun, J. M. U. Vianney, Y. Li, L. Chen, L. Li, F.-Y. Wang, A. Khajepour, and D. Cao, "Proximity based automatic data annotation for autonomous driving," *IEEE/CAA Journal of Automatica Sinica*, vol. 7, no. 2, pp. 395–404, 2020.

[7] C. Sun, Z. Deng, W. Chu, S. Li, and D. Cao, "Acclimatizing the operational design domain for autonomous driving systems," *IEEE Intelligent Transportation Systems Magazine*, 2021.

[8] A. Møgelmose, M. M. Trivedi, and T. B. Moeslund, "Trajectory analysis and prediction for improved pedestrian safety: Integrated framework and evaluations," in *2015 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2015, pp. 330–335.

[9] S. H. Park, B. Kim, C. M. Kang, C. C. Chung, and J. W. Choi, "Sequence-to-sequence prediction of vehicle trajectory via lstm encoder-decoder architecture," in *2018 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2018, pp. 1672–1678.

[10] D. Yang, L. Li, K. Redmill, and Ü. Özgüner, "Top-view trajectories: A pedestrian dataset of vehicle-crowd interaction from controlled experiments and crowded campus," in *2019 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2019, pp. 899–904.

[11] A. Alahi, K. Goel, V. Ramanathan, A. Robicquet, L. Fei-Fei, and S. Savarese, "Social lstm: Human trajectory prediction in crowded spaces," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 961–971.

[12] A. Rudenko, L. Palmieri, M. Herman, K. M. Kitani, D. M. Gavrila, and K. O. Arras, "Human motion trajectory prediction: A survey," *The International Journal of Robotics Research*, vol. 39, no. 8, pp. 895–935, 2020.

[13] O. Styles, A. Ross, and V. Sanchez, "Forecasting pedestrian trajectory with machine-annotated training data," in *2019 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2019, pp. 716–721.

[14] N. Schneider and D. M. Gavrila, "Pedestrian path prediction with recursive bayesian filters: A comparative study," in *German Conference on Pattern Recognition*. Springer, 2013, pp. 174–183.

[15] Y. Li, L. Xin, D. Yu, P. Dai, J. Wang, and S. E. Li, "Pedestrian trajectory prediction with learning-based approaches: A comparative study," in *2019 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2019, pp. 919–926.

[16] A. Rasouli, "Deep learning for vision-based prediction: A survey," *arXiv:2007.00095*, 2020.

[17] A. Gupta, J. Johnson, L. Fei-Fei, S. Savarese, and A. Alahi, "Social gan: Socially acceptable trajectories with generative adversarial networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2255–2264.

[18] X. Li, X. Ying, and M. C. Chuah, "Grip: Graph-based interaction-aware trajectory prediction," in *2019 IEEE Intelligent Transportation Systems Conference (ITSC)*. IEEE, 2019, pp. 3960–3966.

[19] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in neural information processing systems*, 2017, pp. 5998–6008.

[20] F. Giuliari, I. Hasan, M. Cristani, and F. Galasso, "Transformer networks for trajectory forecasting," in *2020 25th International Conference on Pattern Recognition (ICPR)*. IEEE, 2021, pp. 10335–10342.

[21] B. Ivanovic and M. Pavone, "The trajectron: Probabilistic multi-agent trajectory modeling with dynamic spatiotemporal graphs," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 2375–2384.

[22] Y. Ma, X. Zhu, S. Zhang, R. Yang, W. Wang, and D. Manocha, "Trafficpredict: Trajectory prediction for heterogeneous traffic-agents," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 6120–6127.

[23] A. Bhattacharyya, M. Fritz, and B. Schiele, "Long-term on-board prediction of people in traffic scenes under uncertainty," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4194–4202.

[24] C. Yu, X. Ma, J. Ren, H. Zhao, and S. Yi, "Spatio-temporal graph transformer networks for pedestrian trajectory prediction," in *European Conference on Computer Vision*. Springer, 2020, pp. 507–523.

[25] P. Xue, J. Liu, S. Chen, Z. Zhou, Y. Huo, and N. Zheng, "Crossing-road pedestrian trajectory prediction via encoder-decoder lstm," in *2019 IEEE Intelligent Transportation Systems Conference (ITSC)*. IEEE, 2019, pp. 2027–2033.

[26] O. Makansi, O. Cicek, K. Buchicchio, and T. Brox, "Multimodal future localization and emergence prediction for objects in egocentric view with a reachability prior," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 4354–4363.