

# STF: Spatial Temporal Fusion for Trajectory Prediction

Pengqian Han,<sup>1</sup> Partha Roop,<sup>1</sup> Jiamou Liu,<sup>1</sup> Tianzhe Bao,<sup>2</sup> Yifei Wang,<sup>1</sup>

<sup>1</sup> The University of Auckland, Auckland, New Zealand

<sup>2</sup> University of Health and Rehabilitation Sciences, Qingdao, China

phan635@aucklanduni.ac.nz, {p.roop, jiamou.liu}@auckland.ac.nz,

tianzhe.bao@uor.edu.cn, wany107@aucklanduni.ac.nz

**Abstract**—Trajectory prediction is a challenging task that aims to predict the future trajectory of vehicles or pedestrians over a short time horizon based on their historical positions. The main reason is that the trajectory is a kind of complex data, including spatial and temporal information, which is crucial for accurate prediction. Intuitively, the more information the model can capture, the more precise the future trajectory can be predicted. However, previous works based on deep learning methods processed spatial and temporal information separately, leading to inadequate spatial information capture, which means they failed to capture the complete spatial information. Therefore, it is of significance to capture information more fully and effectively on vehicle interactions. In this study, we introduced an integrated 3D graph that incorporates both spatial and temporal edges. Based on this, we proposed the integrated 3D graph, which considers the cross-time interaction information. In specific, we design a Spatial-Temporal Fusion (STF) model including Multi-layer perceptions (MLP) and Graph Attention (GAT) to capture the spatial and temporal information historical trajectories simultaneously on the 3D graph. Our experiment on the ApolloScape Trajectory Datasets shows that the proposed STF outperforms several baseline methods, especially on the long-time-horizon trajectory prediction. Our code is available at <https://github.com/pengqianhan/STF-Spatial-Temporal-Fusion-for-Trajectory-Prediction>

**Index Terms**—Graph Neural Network, Trajectory Prediction, Spatial-Temporal Data Mining

## I. INTRODUCTION

The trajectory prediction is to accurately predict the future path of an object based on its historical position. Trajectory data is a type of spatio-temporal data that contains both spatial and temporal information. Spatial information pertains to the interaction between objects due to their relative position, movement tendency, or relative velocity. Besides physical interaction, there is social interaction between vehicles. Because traffic can not consider all the driving behavior and the drivers do not always follow the traffic rules, sometimes they prefer to reach their destination in the shortest time [1]. Temporal information refers to the time at which the object is at a specific location or moving from one location to another.

The prediction of the trajectories of surrounding vehicles and pedestrians is vital for an autonomous vehicle, enabling it to make informed decisions to avoid collisions. On-board and off-board sensors are used to detect other vehicles, lanes, and crosswalks and generate large amounts of data that can be used

to train high-performance deep learning models. These models can predict feasible and efficient trajectories, improving safety and traffic conditions. Machine learning methods such as Hidden Markov Models (HMMs), Support Vector Machines (SVMs), and Dynamic Bayesian Networks (DBNs) have been used to address this problem. However, these approaches are limited to short time horizons, like 2 seconds, and their prediction performance is unsatisfactory. Deep learning is a subfield of machine learning that utilizes multiple layers in its neural networks and has demonstrated significant efficacy with the increase in data volume [2]. It enables computers to recognize objects from images or translate languages. When applied to trajectory prediction, deep learning methods have proven to outperform conventional machine learning methods. In addition, some work based on deep learning methods like [3]–[7] have demonstrated exceptional performance in long-term trajectory prediction tasks. However, these works process the spatial and temporal separately. As depicted in Fig. 1, certain studies, such as CS LSTM [8], only consider spatial information at the final time step. On the other hand, GRIP++ [9] employs graph neural networks (GNNs) to capture spatial information and subsequently utilize customized convolutional neural networks (CNNs) to capture temporal information. Additionally, Social LSTM [4] and Starnet [10] adopt pooling methods to capture the spatial characteristics of historical trajectories before employing LSTM [11] for capturing temporal information. All of these methods missed the interaction information across the time frames during history time. The vehicles on the road interact with each other, occurring within the same time step and across multiple time steps. Previous research has focused on the interaction between agents. However, they have often treated spatial and temporal information in isolation, resulting in a lack of consideration for interaction related to agents at different time intervals. In this work, we define an integrated 3D graph with spatial edges between the nodes at the same time stamp and temporal edges across the time stamp illustrated in Fig. 2. After defining the 3D graph, all kinds of vehicle interactions are considered, enhancing the model to capture more information from the historical trajectory. Finally, we propose a Spatial-Temporal Fusion (STF) model to predict the future trajectory. Our approach effectively captures more

historical trajectory data information than previous models. We conducted several experiments that demonstrated significantly improved performance over existing methods. The code is available on <https://github.com/>. The main contribution of this work can be summarized as follows:

- This study introduces an integrated 3D graph that incorporates both spatial and temporal edges. The spatial edges represent the spatial interaction between vehicles within the same time frames, while the temporal edges indicate cross-time frame interactions, referred to as spatial-temporal interactions.
- A high-performance trajectory prediction model called STF is proposed in this work. This model effectively captures a wider range of information from historical trajectories, including not only spatial interactions but also spatial-temporal interactions. Additionally, it successfully incorporates temporal information into its predictions.
- Several experiments are conducted to evaluate the performance of the proposed approach, demonstrating superior results compared to previous research papers.

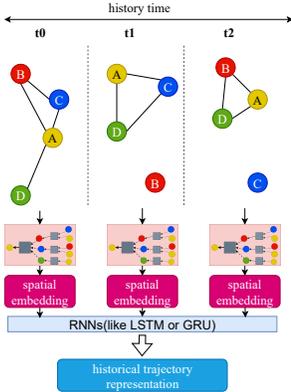


Fig. 1. 2D graph in previous works

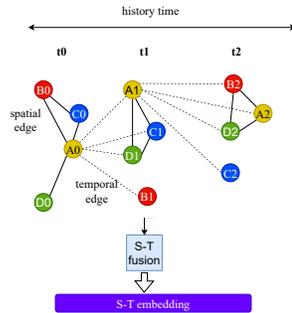


Fig. 2. Integrated 3D graph in S-T Fusion

## II. RELATED WORKS

The methods in previous works are very similar for capturing the temporal information, most of them choosing the RNNs-based models, including LSTM and GRU. However, they apply various algorithms to capture the spatial information of trajectory. These deep learning methods for trajectory prediction can be classified into three categories.

### A. Pooling and CNNs based method

To capture information regarding pedestrian interaction, Social LSTM [4] implemented a grid-based social pooling technique between LSTM cells. Although social pooling addresses the issue of varied numbers of neighboring pedestrians, it does not consider different weights assigned to neighbor agents. TraPHic [12] also utilizes social pooling to capture interaction information. The operation regions in TraPHic’s social pooling consist of a Horizon map and a Neighbor Map, which enable TraPHic to assign different weights to different

agents via social pooling in two maps. It improves the Social LSTM but still misses the cross-time interaction. In CS LSTM [8], shared weights LSTM encoder output is used to make a social tensor, which is put through Convolutional Social Pooling. The Convolutional Social Pooling captures interaction at the last time of the encoder, but it ignores long-range interactions. Social GAN [3] designs a Pooling Module to capture the interaction of all pedestrians in a scene, whereas Social Pooling [4] can only handle people in an  $m \times n$  grid. The relative positions between the target pedestrian and others are fed into an MLP whose output is concatenated with the target’s hidden states. The resulting tensor is then processed by another MLP following a max pooling layer.

### B. GNNs based method

The Social-STGCNN model [5] utilizes Graph Convolutional Networks (GCNs) to capture spatial information and model pedestrian interactions based on a graph defined at each historical frame. Similarly, GRIP++ [9] defines graphs that are then processed by GCNs and 2D temporal convolution layers to capture spatial and temporal information from observed tracks. GSTCN [13] defines a complete graph with weighted edges based on reciprocal distances between nodes to model vehicle interactions. ReCoG [14] uses a directed graph where edges connect vehicles and the target if their distance is below 20 meters. RSBG [15] trains a Recursive Social Behavior Generator to output weighted edges representing latent social relationships, which are then fed into GCN layers. SGCN [16] addresses the issue of superfluous edges by first defining a dense undirected graph and then converting it into a sparse directed graph to model pedestrian interactions. The Social-BiGAT model [17] employs GAT and a latent encoder to process hidden states and generate latent noise, enhancing model multimodality. All the graph-based deep learning methods miss the traffic participants’ interaction cross time because they define the graph at every frame, not the whole historical time.

### C. Attention based method

LaneGCN [18] utilizes attention mechanisms to learn interactions among lanes, actors, and actors themselves. VectorNet [19] employs a hierarchical graph, with subgraphs representing agent trajectories and road features, which are processed using an attention mechanism. HiVT [20] adopts a hierarchical approach by capturing local interactions based on relative positions using an attention mechanism. It then applies a temporal transformer to generate spatial and temporal features of the central agent. Another attention mechanism encodes lane information into keys and values using an MLP and calculates pairwise global interaction. These papers treated temporal information and social interaction separately. The ignorance of cross-time interaction leads to reduced model performance.

## III. PROBLEM FORMULATION

Let  $N$  denote the number of traffic agents,  $T_{his}$  represent the history time, and  $T_{pred}$  indicate the prediction time. The

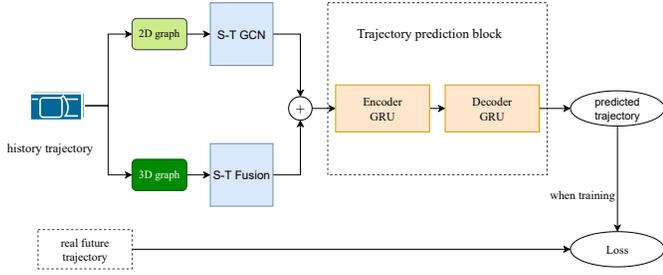


Fig. 3. Scheme of ST fusion model. The 2D graph is defined on all historical time stamps, while the 3D graph is defined based on the entire historical time stamps, including spatial edges and temporal edges

input history trajectories set of agents is denoted by  $\mathbf{X} = [\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^{T_{his}}]$ , where  $\mathbf{x}^t$  represents the coordinates of  $N$  agents at time  $t$ , i.e.,  $\mathbf{x}^t = [x_0^t, y_0^t, x_1^t, y_1^t, \dots, x_{N-1}^t, y_{N-1}^t]$ .

The prediction trajectories are defined as  $\mathbf{Y}$ , which includes future positions from time  $(T_{his} + 1)$  to  $(T_{his} + T_{pred})$ , i.e.,  $\mathbf{Y} = [\mathbf{y}^{T_{his}+1}, \mathbf{y}^{T_{his}+2}, \dots, \mathbf{y}^{T_{his}+T_{pred}}]$ .

#### IV. SPATIAL TEMPORAL FUSION MODEL

##### A. Scheme of model

Fig. 3 illustrates the architecture of the proposed model, which utilizes the historical trajectory of traffic agents, including cars, bikes, and pedestrians, as input. The S-T fusion block extracts both spatial and temporal information simultaneously from this input. Meanwhile, the ST-GCN block separately captures spatial and temporal information. These two blocks' outputs are concatenated before being fed into the Seq2Seq network [21] for predicting future trajectories within a time period of  $T_{pred}$ .

##### B. Integrated 3D graph definition

In order to process the spatial and temporal information simultaneously, this work creates a big graph where each node represents a traffic agent, such as bikes, cars, and pedestrians. An edge is created between two nodes at each time stamp if their distance is shorter than the threshold  $D_{close}$ . These edges are considered **spatial edges**. Additionally, edges between nodes at adjacent time stamps are defined as **temporal edges** used to capture interactions across different time stamps. Fig. 2 displays the spatial and temporal edges. Although other nodes' temporal edges are comparable, they are not depicted in the illustration. Once the historical trajectories are compiled into an integrated 3D graph, an adjacency matrix can be generated, shown in Table I.

##### C. Spatial-Temporal Fusion block

The model takes the location  $(x_i^t, y_i^t)$  of traffic agents during a historical period  $T_{his}$  as input. To better represent this input, two layers of MLP are utilized to increase its dimension. This results in an output dimension of 16 for the MLP, as demonstrated in Equation (1).

$$e_i^t = MLP(x_i^t, y_i^t), \text{ where } e_i^t \in \mathbb{R}^F, F = 16 \quad (1)$$

TABLE I  
ADJACENT MATRIX OF BIGGRAPH

	A0	B0	C0	D0	A1	B1	C1	D1	A2	B2	C2	D2
A0	0	1	1	1	1	1	1	1	0	0	0	0
B0	1	0	1	0	1	1	1	1	0	0	0	0
C0	1	1	0	0	1	1	1	1	0	0	0	0
D0	1	0	0	0	1	1	1	1	0	0	0	0
A1	1	1	1	1	0	0	1	1	1	1	1	1
B1	1	1	1	1	0	0	0	0	1	1	1	1
C1	1	1	1	1	1	0	0	1	1	1	1	1
D1	1	1	1	1	1	0	1	0	1	1	1	1
A2	0	0	0	0	1	1	1	1	0	1	0	1
B2	0	0	0	0	1	1	1	1	1	0	0	1
C2	0	0	0	0	1	1	1	1	0	0	0	0
D2	0	0	0	0	1	1	1	1	1	1	0	0

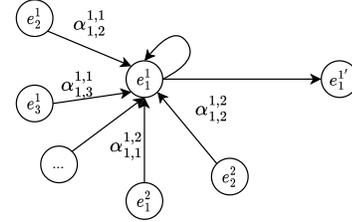


Fig. 4. An illustration of graph attention layer. It allows a node to assign different importance to different nodes within the threshold and in the adjacent time stamps. The new feature of nodes is aggregated from them.

The output of the MLP is provided as input of two layers of GAT, allowing for simultaneous capture of both spatial and temporal information. The input of graph attention layers is Equation (2)

$$h = \left\{ e_1^1, e_2^1, \dots, e_N^1, e_1^2, e_2^2, \dots, e_N^2, \dots, e_1^{T_{his}}, e_2^{T_{his}}, \dots, e_N^{T_{his}} \right\} \quad (2)$$

where  $e_n^t$  is the embedding of  $n$ th node at time  $t$ .

$$\alpha_{u,v}^{i,j} = \frac{\exp(\text{LeakyRelu}(a^T [W e_u^i \oplus W e_v^j]))}{\sum_{k \in \mathbb{N}(u)} \sum_{l \in \{i-1, i, i+1\}} \exp(\text{LeakyRelu}(a^T [W e_u^i \oplus W e_k^l]))} \quad (3)$$

where  $\alpha_{u,v}^{i,j}$  is the attention weights between node  $u$  at time  $i$  and node  $v$  at time  $j$ .  $W \in \mathbb{R}^{F' \times F}$  is a shared trainable matrix for linear mapping.  $F'$  is the output feature dimension, and  $F$  is the input feature dimension.  $\oplus$  denotes the concatenation operation.  $\mathbb{N}(u)$  means the neighbor nodes of  $u$  within a threshold at the same time.  $a \in \mathbb{R}^{2F'}$  is a weight vector.

The updated representation of node  $u$  at time  $i$ , denoted as  $e_u^{i'}$ , is obtained by computing the attention weights of adjacent nodes on target nodes, as specified in Equation (4).

$$e_u^{i'} = \sum_{k \in \mathbb{N}(u)} \sum_{l \in \{i-1, i, i+1\}} \alpha_{u,k}^{i,l} W e_k^l \quad (4)$$

##### D. Spatial-Temporal GCN block

The block records the historical  $x, y$  coordinates of  $N$  traffic agents and processes spatial and temporal information separately. To capture interactions, input data is first fed into a GCN layer, followed by a temporal CNN layer to capture time information. The Spatial-Temporal GCN block consists of three ST-GCN layers, as illustrated in Fig. (6).

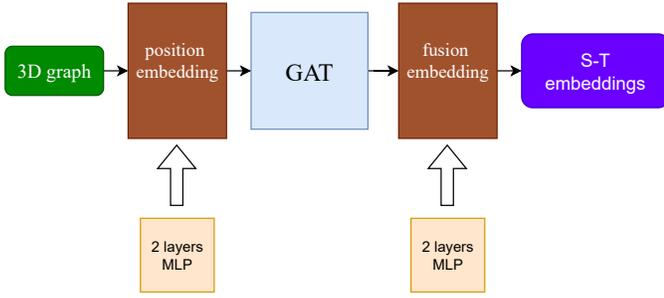


Fig. 5. Spatial-Temporal fusion block

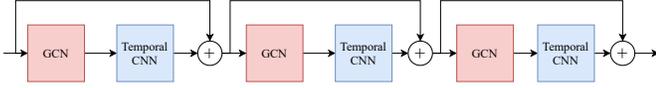


Fig. 6. Spatial-Temporal GCN block

### E. Trajectory prediction block

The trajectory prediction block is a conventional sequence-to-sequence (seq2seq) model comprising GRUs as its components. The seq2seq model excels at processing sequence data and has been widely utilized in prior studies such as [4] and [22]. Previous experiments have demonstrated that the seq2seq model exhibits high performance in trajectory prediction tasks. Furthermore, trajectory prediction belongs to the domain of autoregressive tasks, which are well-suited for LSTMs and GRUs. In this study, we employ the seq2seq model to forecast future trajectories, building upon existing research.

## V. EXPERIMENTS

We trained our model on a Linux server running Ubuntu 22.04.1 LTS with AMD Ryzen Threadripper PRO 3995WX 64-Cores, 512GB RAM, and one NVIDIA GeForce RTX 3090.

### A. Dataset

We train and test our model on the ApolloScape Trajectory Dataset [12], whose framerate are two frames per second. The training sequence lasts for 53 minutes, while the test sequence lasts for 50 minutes. The dataset contains two files, one for training and one for testing. In the training zip file, each file is a 1-minute sequence, and each line in the file contains frame identification, object identification, object type, position x, position y, position z, object length, object width, object height, and heading. The traffic in the dataset is classified into five categories: small vehicles, large vehicles, pedestrians, motorcyclists, bicyclists, and others. The location is based on global coordinates. For the Apollo challenge, the observation time was set to 2 seconds. However, some studies [9] [23] have extended the observation time to 3 seconds, along with a prediction time of 3 seconds. The testing files have the same data structure as the training files, which are used to evaluate the performance of the model.

### B. Metrics

TABLE II  
COMPETITION RESULTS ON APOLLOSCAPE TRAJECTORY DATASET.

Method	WSADE	ADE <sub>v</sub>	ADE <sub>p</sub>	ADE <sub>b</sub>	WSFDE	FDE <sub>v</sub>	FDE <sub>p</sub>	FDE <sub>b</sub>
TrafficPredict [12]	8.5881	7.9467	7.1811	12.8805	24.2262	12.7757	11.1210	22.7912
GRIP++ [9]	2.514	3.948	1.746	3.253	4.026	6.080	2.981	4.913
S2TNet [24]	<b>1.1679</b>	<b>1.9874</b>	<b>0.6834</b>	<b>1.700</b>	2.1798	3.5783	1.3048	3.2151
STF	1.384	2.403	0.799	2.001	<b>1.707</b>	<b>2.972</b>	<b>1.012</b>	<b>2.392</b>

Root Mean Square Error (RMSE): measures the distance between the predictor and the ground truth at time  $t$  [22].

$$RMSE^t = \sqrt{\frac{\sum_{n \in \mathbb{N}} (\hat{y}_n^t - y_n^t)^2}{N}} \quad (5)$$

The Apollo Scape Trajectory dataset provides a weighted sum formula for RMSE, given by Equation (6).

$$RMSE_w = D_v * RMSE_v + D_p * RMSE_p + D_b * RMSE_b \quad (6)$$

where  $D_v = 0.20, D_p = 0.58$ , and  $D_b = 0.22$ . ApolloScape Trajectory Dataset provides these values and points out that they are computed based on the velocity information available in the dataset.  $RMSE_v$  represents the RMSE value for vehicles,  $RMSE_p$  represents the RMSE value for pedestrians, and  $RMSE_b$  represents the RMSE value for bikes.

Average Displacement Error (ADE): calculates the distance between the predicted location of  $N$  agents and ground truth in  $T_{pred}$  time steps, which measures the average prediction performance along the trajectory [8]

$$ADE = \frac{\sum_{t \in \mathbb{T}} \sum_{n \in \mathbb{N}} \|\hat{y}_n^t - y_n^t\|_2}{N \times T_{pred}} \quad (7)$$

where  $\hat{y}_n^t$  is the prediction location of  $n$ th vehicle at time step  $t$ ,  $y_n^t$  is the true location of  $n$ th vehicle at time step  $t$ .  $\mathbb{T} = \{T_{his} + 1, T_{his} + 2, \dots, T_{his} + T_{pred}\}$  and  $\mathbb{N} = \{1, 2, \dots, N\}$ .  $ADE_v, ADE_p$  and  $ADE_b$  are computed separately.

Final Displacement Error (FDE): only computes the distance between the predicted final location and the true final location at the end of final prediction time step  $T_{pred}$  [8]

$$FDE = \frac{\sum_{n \in \mathbb{N}} \|\hat{y}_n^{T_{pred}} - y_n^{T_{pred}}\|_2}{N} \quad (8)$$

where  $\hat{y}_n^{T_{pred}}$  is the prediction location of  $n$ th vehicle at the last prediction time step  $T_{pred}$ ,  $y_n^{T_{pred}}$  is the true location of  $n$ th traffic participant at the last prediction time step  $T_{pred}$ .  $FDE_v, FDE_p$  and  $FDE_b$  are computed separately. The WSADE and WSFDE can be computed through Equation (9) and Equation (10).

$$WSADE = D_c * ADE_v + D_p * ADE_p + D_b * ADE_b \quad (9)$$

$$WSFDE = D_c * FDE_v + D_p * FDE_p + D_b * FDE_b \quad (10)$$

where the  $D_v = 0.20, D_p = 0.58$ , and  $D_b = 0.22$  which is the same as them in Equation (6) provided by the ApolloScape Trajectory Dataset.

TABLE III  
RMSE FOR TRAJECTORY PREDICTION ON APOLLO SCAPE TRAJECTORY

	Object type	RMSE					
		0.5s	1s	1.5s	2s	2.5s	3s
GRIP++	vehicle	<b>1.744</b>	2.577	3.513	4.440	5.335	6.080
	pedestrian	0.515	1.002	1.522	2.022	2.436	2.981
	bike	<b>1.101</b>	2.145	2.924	3.779	4.535	4.913
	weighted RMSE	<b>0.890</b>	1.569	2.228	2.892	3.478	4.026
	all objects	<b>1.384</b>	2.218	3.033	3.919	4.722	5.425
STF	vehicle	1.972	<b>2.0</b>	<b>2.311</b>	<b>2.517</b>	<b>2.649</b>	<b>2.972</b>
	pedestrian	<b>0.50</b>	<b>0.653</b>	<b>0.773</b>	<b>0.887</b>	<b>0.971</b>	<b>1.012</b>
	bike	1.366	<b>1.833</b>	<b>1.925</b>	<b>2.207</b>	<b>2.282</b>	<b>2.392</b>
	weighted RMSE	0.985	<b>1.182</b>	<b>1.334</b>	<b>1.503</b>	<b>1.595</b>	<b>1.707</b>
	all objects	1.586	<b>1.792</b>	<b>2.029</b>	<b>2.265</b>	<b>2.363</b>	<b>2.612</b>

### C. Experiments on the ApolloScape Trajectory Datasets

Table II shows the competition results of TrafficPredict, GRIP++, S2TNet, and STF. The metrics reported include  $WSADE$ ,  $ADE_v$ ,  $ADE_p$ ,  $ADE_b$ ,  $WSFDE$ ,  $FDE_v$ ,  $FDE_p$ , and  $FDE_b$ . Among the listed methods, it is evident that the STF model demonstrates competitive performance across several metrics. For the  $WSADE$  metric, STF achieves a value of 1.384, which is lower than TrafficPredict (8.5881) and GRIP++ (2.514), indicating improved accuracy. Similarly, for the  $ADE_v$ ,  $ADE_p$ , and  $ADE_b$  metrics, the STF model outperforms TrafficPredict and GRIP++, showing lower values and better accuracy in trajectory prediction. But the higher values than the result of S2TNet, which means the S2TNet has better performance than STF in short-time prediction. Moving on to the  $WSFDE$ ,  $FDE_v$ ,  $FDE_p$ , and  $FDE_b$  metrics, the STF model excels, demonstrating the lowest values among the listed methods. Specifically, STF achieves an impressive  $WSFDE$  of 1.707, an  $FDE_v$  of 2.972, an  $FDE_p$  of 1.012, and an  $FDE_b$  of 2.392. These results indicate that the STF model can predict long-term trajectories with greater precision and accuracy compared to TrafficPredict, GRIP++, and S2TNet.

In summary, the STF model stands out as the best-performing method in trajectory prediction based on the competition results. It surpasses TrafficPredict, GRIP++, and S2TNet in terms of  $FDE$  metrics, demonstrating superior accuracy and precision in trajectory prediction on long-time trajectory prediction.

Due to the lack of  $RMSE$  values at every frame in the TrafficPredict and S2TNet models, we evaluate our model by comparing it with GRIP++ based on the  $RMSE$  at each prediction frame. Furthermore, our model incorporates GRIP++ as its backbone. The outcomes derived from GRIP++ can be regarded as an ablation experiment. The experiment results, as shown in Table III, compare the  $RMSE$  for trajectory prediction on the Apollo Scape Trajectory dataset between the GRIP++ model and the STF. The  $RMSE$  values are reported for different prediction time intervals, ranging from 0.5s to 3s.

Overall, the STF outperforms the GRIP++ model in terms of trajectory prediction accuracy. Specifically, the STF achieves lower  $RMSE$  values for most object types and prediction time intervals compared to GRIP++. For the car object type, the STF obtains slightly higher  $RMSE$  values at shorter prediction time intervals (0.5s) but achieves better performance

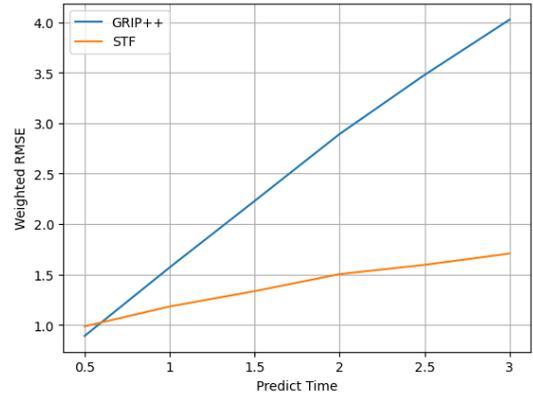


Fig. 7. Weighted RMSE in prediction time

at longer prediction time intervals (1s, 1.5s, 2s, 2.5s, and 3s). For pedestrians, the STF consistently exhibits lower  $RMSE$  values across all prediction time intervals, especially for long-term prediction. For the bike object type, the STF achieves better results at longer prediction time intervals (1s, 1.5s, 2s, 2.5s, and 3s), while the  $RMSE$  values are comparable or slightly higher at shorter intervals (0.5s). For the weighted sum of these three kinds of objects computed according to Equation (6), the STF shows much better on the long-term (1s, 1.5s, 2s, 2.5s, and 3s) trajectory prediction. The STF model demonstrates superior performance compared to the GRIP++ model in terms of all object types. This includes vehicles, pedestrians, bikes, and other objects. The STF consistently achieves lower  $RMSE$  values across long prediction time intervals (1s, 1.5s, 2s, 2.5s, and 3s). Fig. 7 displays the plot of the Weighted  $RMSE$  for GRIP++ and STF at various predicted times. This figure highlights the superior performance of STF in accurately predicting future trajectories over an extended period.

## VI. CONCLUSION

This paper presents the STF model, designed to capture spatial and temporal information from historical trajectory data simultaneously. The model addresses the information missing problem in previous works while mitigating the issue of error accumulation. Experiments are conducted on the ApolloScape Trajectory Datasets. The results demonstrate superior performance compared to previous models. Our model exhibits a significant advantage, particularly in long-term trajectory prediction performance. Accurate long-term predictions ensure that autonomous vehicles and robotics can proactively avoid potential collisions, providing them with more time to respond, ultimately resulting in enhanced safety.

## REFERENCES

- [1] W. Wang, L. Wang, C. Zhang, C. Liu, L. Sun *et al.*, "Social interactions for autonomous driving: A review and perspectives," *Foundations and Trends® in Robotics*, vol. 10, no. 3-4, pp. 198–376, 2022.
- [2] A. Zhang, Z. C. Lipton, M. Li, and A. J. Smola, "Dive into deep learning," *arXiv preprint arXiv:2106.11342*, 2021.

- [3] 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.
- [4] A. Alahi, K. Goel, V. Ramanathan, A. Robicquet, L. Fei-Fei, and S. Savarese, "Social LSTM: Human Trajectory Prediction in Crowded Spaces," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. Las Vegas, NV, USA: IEEE, Jun. 2016, pp. 961–971.
- [5] A. Mohamed, K. Qian, M. Elhoseiny, and C. Claudel, "Social-stgcn: A social spatio-temporal graph convolutional neural network for human trajectory prediction," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 14 424–14 432.
- [6] C. Yu, X. Ma, J. Ren, H. Zhao, and S. Yi, "Spatio-temporal graph transformer networks for pedestrian trajectory prediction," in *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XII 16*. Springer, 2020, pp. 507–523.
- [7] M. Mendieta and H. Tabkhi, "Carpe posterum: A convolutional approach for real-time pedestrian path prediction." *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 3, p. 2346–2354, Sep 2022. [Online]. Available: <http://dx.doi.org/10.1609/aaai.v35i3.16335>
- [8] N. Deo and M. M. Trivedi, "Convolutional Social Pooling for Vehicle Trajectory Prediction," May 2018.
- [9] X. Li, X. Ying, and M. C. Chuah, "Grip++: Enhanced graph-based interaction-aware trajectory prediction for autonomous driving," *arXiv preprint arXiv:1907.07792*, 2019.
- [10] Y. Zhu, D. Qian, D. Ren, and H. Xia, "Starnet: Pedestrian trajectory prediction using deep neural network in star topology," in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2019, pp. 8075–8080.
- [11] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [12] 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.
- [13] Z. Sheng, Y. Xu, S. Xue, and D. Li, "Graph-based spatial-temporal convolutional network for vehicle trajectory prediction in autonomous driving," *IEEE Transactions on Intelligent Transportation Systems*, 2022.
- [14] X. Mo, Y. King, and C. Lv, "Recog: A deep learning framework with heterogeneous graph for interaction-aware trajectory prediction," *arXiv preprint arXiv:2012.05032*, 2020.
- [15] J. Sun, Q. Jiang, and C. Lu, "Recursive social behavior graph for trajectory prediction," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 660–669.
- [16] L. Shi, L. Wang, C. Long, S. Zhou, M. Zhou, Z. Niu, and G. Hua, "Sgc: Sparse graph convolution network for pedestrian trajectory prediction," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 8994–9003.
- [17] V. Kosaraju, A. Sadeghian, R. Martín-Martín, I. Reid, H. Rezatofighi, and S. Savarese, "Social-bigat: Multimodal trajectory forecasting using bicycle-gan and graph attention networks," *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [18] M. Liang, B. Yang, R. Hu, Y. Chen, R. Liao, S. Feng, and R. Urtasun, "Learning lane graph representations for motion forecasting," in *European Conference on Computer Vision*. Springer, 2020, pp. 541–556.
- [19] J. Gao, C. Sun, H. Zhao, Y. Shen, D. Anguelov, C. Li, and C. Schmid, "Vectornet: Encoding hd maps and agent dynamics from vectorized representation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 11 525–11 533.
- [20] Z. Zhou, L. Ye, J. Wang, K. Wu, and K. Lu, "Hivt: Hierarchical vector transformer for multi-agent motion prediction," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 8823–8833.
- [21] I. Sutskever, O. Vinyals, and Q. V. Le, "Sequence to sequence learning with neural networks," *Advances in neural information processing systems*, vol. 27, 2014.
- [22] H. Song, W. Ding, Y. Chen, S. Shen, M. Y. Wang, and Q. Chen, "Pip: Planning-informed trajectory prediction for autonomous driving," in *European Conference on Computer Vision*. Springer, 2020, pp. 598–614.
- [23] 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.
- [24] W. Chen, F. Wang, and H. Sun, "S2tnet: Spatio-temporal transformer networks for trajectory prediction in autonomous driving," in *Asian Conference on Machine Learning*. PMLR, 2021, pp. 454–469.