

TrajMamba: An Ego-Motion-Guided Mamba Model for Pedestrian Trajectory Prediction from an Egocentric Perspective

Yusheng Peng¹, Gaofeng Zhang² and Liping Zheng^{2*}

Abstract—Future trajectory prediction of a tracked pedestrian from an egocentric perspective is a key task in areas such as autonomous driving and robot navigation. The challenge of this task lies in the complex dynamic relative motion between the ego-camera and the tracked pedestrian. To address this challenge, we propose an ego-motion-guided trajectory prediction network based on the Mamba model. Firstly, two Mamba models are used as encoders to extract pedestrian motion and ego-motion features from pedestrian movement and ego-vehicle movement, respectively. Then, an ego-motion-guided Mamba decoder that explicitly models the relative motion between the pedestrian and the vehicle by integrating pedestrian motion features as historical context with ego-motion features as guiding cues to capture decoded features. Finally, the future trajectory is generated from the decoded features corresponding to the future timestamps. Extensive experiments demonstrate the effectiveness of the proposed model, which achieves state-of-the-art performance on the PIE and JAAD datasets.

I. INTRODUCTION

Predicting future trajectories of pedestrians is a key aspect of pedestrian behavior analysis, which is essential for ensuring the safe driving of autonomous vehicles and the secure navigation of mobile robots in dynamic crowd environments. Most research [1], [2] model and analyze pedestrian movement from a bird’s-eye view to predict the future trajectories of pedestrians. This type of method is only suitable for third-person perspective videos captured by surveillance cameras or drones. Autonomous vehicles and mobile robots typically perceive their surroundings from a first-person perspective. Specifically, these ego cameras usually move in tandem with the terminal agent. Nevertheless, the dual motion of the ego camera and the tracked pedestrian poses significant challenges in predicting the pedestrian’s future trajectory from the egocentric perspective.

The core objective of first-person perspective pedestrian trajectory prediction is to accurately forecast the future bounding box sequence of the target pedestrian within the ego-camera coordinate system. The primary challenge stems from the fact that the apparent motion observed in the image

is the result of the superimposed projection of both the pedestrian’s intrinsic motion and the carrier’s (vehicle or camera) motion onto the two-dimensional imaging plane. Consequently, the prediction accuracy heavily relies on the effective modeling of the relative motion between these two sources. Most existing approaches tackle this challenge through feature-level fusion strategies, where pedestrian motion and vehicle ego-motion are combined in the latent space using early or late fusion techniques, and subsequently fed into a decoder for end-to-end prediction. For example, methods [3], [4] employ LSTM networks as encoders to extract pedestrian motion features and vehicle ego-motion features separately, and then concatenate these two feature types as input to the decoder LSTM. Similarly, method [5] adopts Transformer to replace LSTM for both encoding and decoding. Meanwhile, method [6] concatenates the two types of features at the embedding level and processes them through a Transformer to obtain fused encoded representations, which are then decoded to generate future trajectories. In all these approaches, the two types of motion features are integrated into a unified representation during the encoding stage, leaving the decoder to generate future predictions. This technical paradigm results in a rather ambiguous modeling of the relative motion relationship between the two motion sources, making it difficult to explicitly analyze how ego-motion dynamically regulates pedestrian movement.

To address these limitations, we propose TrajMamba, a novel ego-motion-guided framework for egocentric pedestrian trajectory prediction. Specifically, we first employ two Mamba encoders to extract pedestrian motion features and ego-motion features, respectively. An ego-motion-guided Mamba decoder is then designed, which takes pedestrian motion features as historical context and ego-motion features as future guiding cues. By leveraging Mamba’s sequential modeling capacity, the decoder explicitly captures the dynamic modulation mechanism of ego-motion on pedestrian movement to infer future motion features. Finally, these decoded features are mapped to future trajectories via a prediction head. The main contributions of this work are as follows:

- We propose TrajMamba, a novel ego-motion-guided framework based on the Mamba architecture, for predicting future pedestrian trajectories from an egocentric perspective.
- We design an ego-motion-guided Mamba decoder that explicitly models the relative motion between the pedestrian and the camera by integrating pedestrian motion

*This work was supported in part by the National Natural Science Foundation of China under Grant 62372152, in part by the Fundamental Research Funds for the Central Universities of China under Grant JZ2023HGQB0481, and in part by China Postdoctoral Science Foundation under Grant 2023M740961

¹Yusheng Peng is with the School of Computer Science and Information Engineering, Hefei University of Technology, Hefei 230601, China (e-mail: wisionpeng@hfut.edu.cn).

²Gaofeng Zhang and Liping Zheng are with the School of Software, Hefei University of Technology, Hefei 230601, China (e-mail: g.zhang@hfut.edu.cn, zhenglp@hfut.edu.cn).

*Corresponding author

features as historical context with ego-motion features as guiding cues for future prediction.

- Extensive experiments on the PIE and JAAD datasets demonstrate that TrajMamba outperforms baseline methods on the ADE, FDE, ARB, and FRB metrics.

II. RELATED WORK

A. Pedestrian Trajectory Prediction

Current methods for predicting pedestrian trajectories are classified into those based on the bird’s eye perspective and those based on the egocentric perspective. Bird’s-eye perspective-based methods focus on predicting the 2D coordinates of pedestrians’ future locations in bird’s eye videos captured by surveillance cameras or drones. As a typical sequence prediction task, many works [7], [8] employ the recurrent neural networks (RNNs) or its variants (LSTM and GRU) to capture temporal dependencies of pedestrian’s movement, and utilize pooling mechanism [9] or graph neural networks [10] to model spatial interactions among pedestrians. Recent works [11], [12] have utilized Transformers to replace networks such as LSTM, leveraging their attention mechanisms to model temporal dependencies of pedestrian movement and spatial interactions. Additionally, some researchers [13], [14] integrate recently popular diffusion models and introduce new architectures to predict pedestrian trajectories. Moreover, some works [15], [16] leverage the visual information to improve the prediction performance.

On the other hand, egocentric perspective-based methods focus on predicting future bounding boxes of tracked pedestrians from egocentric videos captured by vehicle-mounted or wearable cameras. Inspired by bird’s-eye view future localization methods, some researchers [17], [18], [19] use LSTM or GRU to process the bounding boxes of tracked pedestrians to model their motion dynamics and predict the bounding boxes of future locations. Furthermore, the researchers [20], [21] thoroughly investigate pedestrian behavior attributes, incorporating pose or behavior labels as auxiliary information in their models, which enhances the accuracy of future localization. Moreover, some methods [22], [5] also incorporate scene semantics to more accurately predict the future locations of pedestrians in diverse environments. The egocentric view is unique in that the camera itself is in motion, presenting new challenges for predicting a pedestrian’s future location. To tackle this, researchers [23], [4] extract the camera’s ego-motion features using the vehicle’s travel data or the camera’s IMU signals. By combining the pedestrian’s motion with the camera’s ego-motion, they enhance the accuracy of future localization predictions. Different from the above works, we innovatively introduce the Mamba model for pedestrian motion modeling and ego-motion modeling, which improves the efficiency while maintaining the modeling ability.

B. State Space Models

Recently, state-space models (SSMs) [24], [25] have garnered significant attention due to their excellent performance

in sequence modeling and inference. Mamba [26] introduces a selective state space model architecture that integrates time-varying parameters into the SSM framework and proposes a hardware-aware algorithm to enhance the efficiency of training and inference processes. Furthermore, Mamba-2 [27] refines this by linking SSMs to attention variants, achieving 2-8x speedup and performance comparable to transformers. SSMs and Mamba models are widely and successfully used in many tasks. For instance, Zhu et al. [28] introduce a new generic vision backbone called Vision Mamba (Vim), which incorporates bidirectional Mamba blocks. Vim leverages position embeddings to mark image sequences and compresses visual representations using bidirectional state-space models. Similarly, Park et al. [29] introduce VideoMamba for video analysis, which tackles the unique challenge of integrating non-sequential spatial with sequential temporal information in video processing through spatio-temporal forward and backward SSM. Tang et al. [30] propose a novel spatio-temporal graph Mamba (STG-Mamba) for the music-guided dance video synthesis task, i.e., to translate the input music to a dance video. He et al. [31] propose a decomposed spatio-temporal Mamba (DST-Mamba) for traffic prediction. Zhang et al. [32] introduce a Mamba-based point cloud network named Point Cloud Mamba, which incorporates several novel techniques to help Mamba better model point cloud data. Inspired by these, we explore a novel application of the Mamba model for pedestrian trajectory prediction from an egocentric perspective.

III. METHOD

A. Overview

Our goal is to forecast the future bounding boxes of pedestrians from an egocentric perspective. In this section, we propose a novel ego-motion-guided Mamba model named TrajMamba. As shown in Fig. 1, the proposed TrajMamba model includes four modules: pedestrian motion encoder (PME), ego-motion encoder (EME), ego-motion guide decoder (EMGD), and future trajectory generator (FTG) module. Firstly, the pedestrian motion features and the ego-motion features are extracted through the PME and EME modules, respectively. Then, the pedestrian’s motion features and ego-motion features are sent to the EMGD module to output the decoded features by modeling the relative motion relationship between the pedestrian and the ego-vehicle. Finally, the decoded features of future timestamps are then used to generate future trajectory through the FTG module.

B. Problem Formulation

Suppose that the past bounding boxes of the tracked pedestrian at current frame t are denoted as a sequence $B_{past} = (B_{T-m}, \dots, B_T)$, and its future bounding boxes are denoted as a sequence $B_{future} = (B_T, \dots, B_{T+n})$, where $B_t = (x_t, y_t, w_t, h_t)$, and m and n are the lengths of past and future timestamps. The movement information of the ego-vehicle is represented as $V_{past} = (v_{T-m}, \dots, v_T)$. Our goal is to predict the future bounding boxes $B_{future} = (\tilde{B}_T, \dots, \tilde{B}_{T+n})$ of the tracked pedestrian from its past

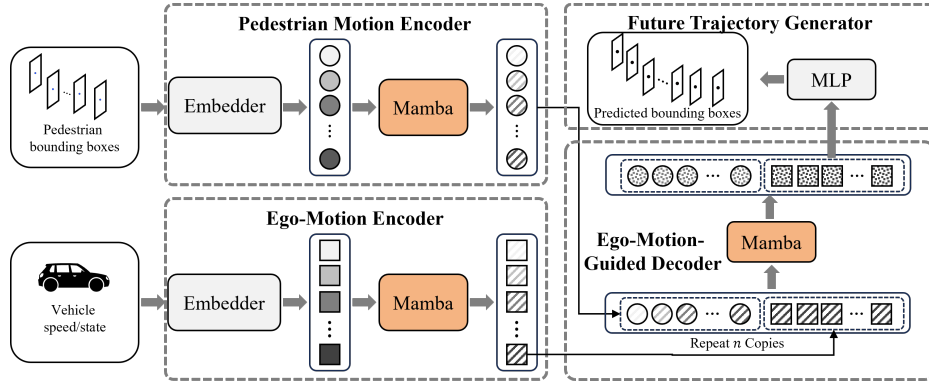


Fig. 1: Overview of the proposed TrajMamba.

bounding boxes B_{past} and the ego-vehicle's past movement information V_{past} .

C. Input and Output Representation

Existing approaches to pedestrian motion representation can be categorized into two paradigms: one utilizes the center point coordinates with bounding box width and height [19], [33], while the other employs the top-left and bottom-right corner coordinates [5], [34]. Both paradigms typically maintain consistent input-output representations. Recently, differentiated representation strategies have emerged. For example, method [22] inputs corner coordinates and velocities but outputs regressed corner-based bounding boxes. Method [35] inputs center coordinates with velocity and acceleration, yet outputs only the predicted center coordinates. Motivated by this, we propose a differentiated motion representation method. Specifically, given a pedestrian's historical bounding box sequence, we compute the center point displacement (velocity) and scale variation between adjacent frames as follows:

$$v_t^x = x_t - x_{t-1}, v_t^y = y_t - y_{t-1} \quad (1)$$

$$\Delta w = w_t - w_{t-1}, \Delta h = h_t - h_{t-1} \quad (2)$$

Based on this, the motion state of the pedestrian at time step t is constructed as $M_t = (x_t, y_t, w_t, h_t, v_t^x, v_t^y, \Delta w, \Delta h)$.

On the other hand, inspired by the constant velocity assumption in [36], we introduce the Constant Velocity and Constant Scaling (CV-CS) assumption. The key innovation lies in predicting residuals relative to a physically motivated reference trajectory, which not only eases network optimization but also injects valuable prior knowledge into the learning process. Specifically, based on the CV-CS assumption, pedestrian motion approximately maintains constant velocity and constant scale over short time intervals. Therefore, we compute the average velocity (\bar{v}^x, \bar{v}^y) and average scale variation rate (\bar{h}, \bar{w}) using the last five frames of historical observations. According to this assumption, the reference bounding box at any future time step $t + \tau$ can be inferred:

$$\hat{x}_{t+\tau} = x_T + \bar{v}^x \cdot \tau, \hat{y}_{t+\tau} = y_T + \bar{v}^y \cdot \tau \quad (3)$$

$$\hat{w}_{t+\tau} = w_T \cdot (1 + \bar{w})^\tau, \hat{h}_{t+\tau} = h_T \cdot (1 + \bar{h})^\tau \quad (4)$$

We define the residual offset $\Delta_{T+\tau}$ of the ground-truth future bounding box $B_{T+\tau}$ relative to the reference $\hat{B}_{T+\tau}$. The network is tasked with predicting this offset, denoted as $\hat{\Delta}_{T+\tau}$. Finally, the predicted bounding box $\hat{B}_{T+\tau}$ is obtained by adding the predicted offset $\hat{\Delta}_{T+\tau}$ to the reference $\hat{B}_{T+\tau}$.

This method reformulates the prediction target into a residual form, which not only preserves the physical prior of the CV-CS assumption but also allows the network to flexibly correct deviations from actual motion. Consequently, it achieves more accurate trajectory prediction in complex dynamic scenarios.

D. Pedestrian Motion Encoder

To precisely predict the future trajectory of the tracked pedestrian, it is of vital importance to initially model and analyze their past movement dynamics and patterns so as to capture the underlying principles of their behavior. For a pedestrian's past bounding boxes B_{past} , we begin by calculating the motion state $M_t = (x_t, y_t, w_t, h_t, v_t^x, v_t^y, \Delta w, \Delta h)$. Then, we utilize a multilayer perceptron (MLP) as the embedder layer to transform the motion state sequence $\{M_{T-m}, \dots, M_T\}$ into embedding features $E_{pm} \in \mathbb{R}^{(T-m) \times D}$. The D is the dimension of each embedding feature. Existing methods always employ LSTM [37], GRU [38], or Transformer [5] model as the encoder to capture motion features. Different from the previous methods, we use the Mamba model as the encoder, aiming to effectively extract motion features and improve efficiency by taking advantage of its linear time series modeling ability and efficient long sequence processing performance. The encode operation of the Mamba model is formulated as follows:

$$F_{pm} = Mamba(E_{pm}) \quad (5)$$

E. Ego-Motion Encoder

Modeling the motion of the ego camera is key to understanding the motion dynamics of the tracked pedestrian in an egocentric perspective. Existing works [39], [40], [41] utilize CNNs to implicitly capture ego-motion features from image streams. Other works usually use LSTMs [17], [3] or Transformers [42], [5] to model the movement information of ego vehicles to capture ego motion features. Inspired by this, we use a Mamba model to extract the motion

features of ego vehicles. Specifically, given the speed or state information v_t of the vehicle at t timestamp, an MLP layer is used as an embedder to transform this information into an embedded feature e_t and obtain the embedded feature sequence $E_{em} = \{e_{t-m}, \dots, e_t\}$ of past m timestamps. Then, a Mamba model is used as an encoder to transform the embedded features E_{em} into ego-motion features F_{em} through sequential modeling. The encode operation of the Mamba model is formulated as follows:

$$F_{em} = \text{Mamba}(E_{em}) \quad (6)$$

F. Ego-Motion-Guided Decoder

Existing methods primarily integrate pedestrian motion and vehicle ego-motion into a unified encoded feature through early [6], late [3], [18], or cross [44], [45] fusion strategies, and then generate future trajectories by decoding this fused representation. Although such schemes achieve collaborative modeling of the two motion sources, they struggle to explicitly analyze the dynamic modulation mechanism of ego-motion on pedestrian movement. To overcome the limitation of implicit fusion, we introduce an explicit modeling framework. In this framework, pedestrian motion features are treated as historical context, and ego-motion features are leveraged as explicit conditional cues for future prediction. By employing sequential modeling within the decoder, the network explicitly learns how ego-motion dynamically modulates pedestrian movement. This allows for accurate inference of future pedestrian states from the egocentric viewpoint.

Given the superior performance of the Mamba model over GRU, LSTM, and Transformer networks in various sequence prediction tasks, we adopt it as our decoder. In our design, pedestrian motion features serve as the historical context for the decoder, while the vehicle's motion features at the last observed timestamp (T) are used as explicit conditional guidance for the future prediction horizon. By leveraging the structured modeling mechanism of the Mamba network, our approach explicitly learns how ego-motion dynamically modulates pedestrian movement. This enables effective inference of future pedestrian states from an egocentric perspective, ultimately outputting the decoded features for future time steps:

$$F_{de} = \text{Mamba}(F_{in}) \quad (7)$$

where the input feature sequence F_{in} is composed of features sequence F_m and T_{pred} copies of the ego-motion feature at timestamp T_{pred} . Through the decoding operation, a set of decoded features F_{de} that represent the future motion of the pedestrian is captured.

G. Future Trajectory Generator

The decoder Mamba model captures a set of decoded features $F_{de}^{t:t+n}$ that represent the future motion of a pedestrian. Next, an MLP acts as a generator to map the decoded features of future timestamps into a set of future trajectories. The future trajectory generator is formulated as follows:

$$\hat{F}_{future} = \text{MLP}(F_{de}^{t:t+n}) \quad (8)$$

A. Datasets and Metrics

We conduct experiments to train and evaluate the prediction performance of the proposed model on the widely used benchmark datasets, namely JAAD [46] and PIE [43]. JAAD is a dataset for studying joint attention in the context of autonomous driving, which is widely used in trajectory prediction from egocentric views and crossing intention prediction research. To this end, JAAD dataset provides a richly annotated collection of 346 short video clips (5-10 sec long) extracted from over 240 hours of driving footage. PIE is a dataset for studying pedestrian behavior in traffic, which contains over 6 hours of footage recorded in typical traffic scenes with on-board camera. It also provides accurate vehicle information from OBD sensor (vehicle speed, heading direction and GPS coordinates) synchronized with video footage. Following the benchmark protocol, the datasets are randomly split into training (50%), validation (10%), and test (40%) sets.

Similar to existing works [3], [22], [5], we assess the performance of trajectory prediction by using four metrics: Average Displacement Error (ADE), Final Displacement Error (FDE), Average Root Mean Square Error (ARB), and Final Root Mean Square Error (FRB). ADE and FDE measure the displacement error of center point coordinates between the predicted bounding boxes and the ground truth bounding boxes, and ARB and FRB measure the mean square error between the predicted bounding boxes and the coordinates of the groundtruth bounding boxes.

B. Experimental Setup

We train and evaluate the TrajMamba model on an Nvidia GTX 4080 GPU with CUDA 11.6 and PyTorch 2.0.0. The dimensions of pedestrian motion features and ego-motion features are set to 256. The past and future timestamps are set to 15 (0.5s) and 45 (1.5s). We use the Adam optimizer with an initial learning rate of 0.0001 and smooth L_1 loss as the loss function. In particular, for the PIE dataset, the vehicle speed is used to represent the vehicle's movement information to capture ego-motion features, while in the JAAD dataset, the vehicle behaviors (0: stopped, 1: moving slow, 2: moving fast, 3: decelerating, and 4: accelerating) are used as the movement information.

C. Quantitative Evaluation

We evaluated the proposed TrajMamba model against eight existing methods, including FOL [40], FPL [39], B-LSTM [17], and two variations of the methods in the PIE dataset [43], BiPed [3], PedFormer[22], and two general models BiTrap [37] and SGNet [38], PEvT [41], MTN [42], and the model proposed by Zhang et al. [5].

The experimental results between TrajMamba and baseline methods on the PIE and JAAD datasets are presented in Table I. For 1s prediction (upper part), TrajMamba consistently outperforms all methods. Notably, it surpasses the previous best model PedFormer by 40.98%/ 44.02%/ 24.03%/ 32.91% on ADE/ FDE/ ARB/ FRB on PIE, and achieves

TABLE I: Performance of the proposed TrajMamba and other existing models on the PIE and JAAD datasets.

Models	Publications	Years	PIE				JAAD			
			ADE ↓	FDE ↓	ARB ↓	FRB ↓	ADE ↓	FDE ↓	ARB ↓	FRB ↓
			Observation frames: 15 (0.5s)				Prediction frames: 30 (1.0s)			
FOL [40]	ICRA	2019	73.87	164.53	78.16	143.49	61.39	126.97	70.12	129.17
FPL [39]	CVPR	2018	56.66	132.23	/	/	42.24	86.13	/	/
B-LSTM [17]	CVPR	2018	27.09	66.74	37.41	75.87	28.36	70.22	39.14	79.66
PIE_traj [43]	ICCV	2019	21.82	53.63	27.16	55.39	23.49	50.18	30.40	57.17
PIE_full [43]	ICCV	2019	19.50	45.27	24.40	49.09	22.83	49.44	29.52	55.43
BiPed [3]	ICCV	2021	15.21	35.03	19.62	39.12	20.58	46.85	27.98	55.07
PedFormer [22]	ICRA	2023	13.08	30.35	15.27	32.79	17.89	41.63	24.56	48.82
BiTrap [37]	RAL	2021	8.83	18.52	12.61	22.97	19.02	39.28	22.57	40.88
SGNet [38]	RAL	2022	8.35	17.83	12.32	22.58	14.90	30.88	19.19	34.80
Ours	/	/	7.72	16.99	11.60	22.00	13.90	30.76	18.04	34.79
			Observation frames: 15 (0.5s)				Prediction frames: 45 (1.5s)			
PEvT [41]	CVPR	2021	19.15	45.98	/	/	21.08	49.08	/	/
MTN [42]	IJCAI	2021	18.89	45.50	/	/	20.90	48.55	/	/
Zhang et al. [5]	TIP	2024	17.41	44.92	/	/	17.92	41.33	/	/
Ours	/	/	13.01	31.08	20.63	41.79	23.39	55.11	32.06	63.65

TABLE II: Ablation results on input and output representation.

Input	Output	ADE	FDE	ARB	FRB
		PIE			
LT-BR	LT-BR	8.72	18.69	24.68	27.88
xywh	xywh	13.64	24.84	51.50	78.76
LT-BR + v	LT-BR	8.58	18.20	13.15	24.16
xywh+v+s	CV-CS-Offset	7.72	16.99	11.60	22.00
		JAAD			
LT-BR	LT-BR	17.27	35.17	23.24	41.95
xywh	xywh	26.66	42.12	52.21	79.19
LT-BR + v	LT-BR	15.93	32.76	20.74	38.15
xywh+v+s	CV-CS-Offset	13.90	30.76	18.04	34.79

22.30%/ 26.11%/ 26.55%/ 28.74% improvements on JAAD. Against general approaches BiTrap and SGNet, TrajMamba yields gains of 12.57%/ 8.26%/ 8.01%/ 4.22% and 7.54%/ 4.71%/ 5.84%/ 2.57% on PIE, and 26.92%/ 21.69%/ 20.07%/ 14.90% and 6.71%/ 0.39%/ 5.99%/ 0.03% on JAAD, respectively. For 1.5s prediction (lower part), TrajMamba achieves best performance on PIE but lags behind PEvT, MTN, and Zhang et al.’s method on JAAD. This discrepancy likely stems from differences in scene characteristics and vehicle motion representations (speed vs. behavior labels) between datasets.

D. Ablation Study

The ablation study on input-output representation:

To evaluate the effectiveness of the proposed input-output representation, we conducted an ablation study, with results reported in Table II. In the table, LT-BR denotes the coordinates of the top-left and bottom-right corners of the bounding box; xywh represents the center coordinates along with the width and height; v indicates the velocity along x and y directions; s refers to the scale variation on weight and height; and CV-CS-Offset corresponds to the residual offset derived

from the constant velocity and constant scale assumption. Experiment results demonstrate that incorporating velocity information significantly improves prediction accuracy, highlighting the critical role of encoding motion dynamics. Furthermore, the representation that jointly encodes position, velocity, and scale variation—combined with the offset output derived from the CV-CS assumption—further enhances trajectory smoothness and physical plausibility, providing a more robust and interpretable motion feature foundation for the model.

The effect of ego-motion-guided decoder: To verify the effectiveness of the proposed ego-motion-guided decoder, we conducted multiple sets of comparative experiments, with results summarized in Table III. For each dataset, we designed two categories of ablation experiments. The first category follows the classic architecture: pedestrian motion features and ego-motion features are extracted separately by encoders, then fused and processed by a decoder to generate future trajectories. The second category corresponds to our proposed architecture, where pedestrian motion features and ego-motion features are extracted separately and then fed into the designed ego-motion-guided decoder to generate future trajectories. Specifically, within each category, we employ LSTM, GRU, Transformer, and Mamba as the encoder/decoder backbones for comprehensive comparison.

From the results presented in the table, it can be observed that for the same base model (with the exception of LSTM), the variant constructed with our proposed ego-motion-guided decoding mechanism (EMGD) consistently achieves more accurate predictions than the classic post-fusion decoding mechanism (PFD). This improvement arises because traditional methods provide a fused motion representation that mixes pedestrian motion with ego-motion during decoding, leading to ambiguity and making it difficult for the model to disentangle the relative motion relationship between the pedestrian and the vehicle. In contrast, our proposed ego-motion-guided decoder takes pedestrian motion features as

TABLE III: Ablation Results on Ego-motion-guided decoder. (PFD: post-fusion decoding mechanism, EMGD: ego-motion-guided decoding mechanism)

	Model	ADE	FDE	ARB	FRB	Flops	Params
PIE							
PFD	LSTM	7.98	17.82	12.14	23.29	128.27G	3.16M
	GRU	7.97	17.98	12.25	23.72	100.68G	2.51M
	Transformer	7.94	17.86	12.25	23.52	87.48G	2.31M
	Mamba	7.91	17.87	11.99	23.19	30.11G	1.32M
EMGD	LSTM	8.14	18.26	12.34	23.57	63.64G	1.78M
	GRU	7.97	17.87	12.13	23.28	49.34G	1.39M
	Transformer	7.72	17.13	11.84	22.51	6.35G	0.20M
	Mamba	7.72	16.99	11.6	22	6.35G	0.20M
JAAD							
PFD	LSTM	15.66	35.37	20.16	39.07	220.99G	3.16M
	GRU	15.19	34.51	19.71	38.29	173.46G	2.51M
	Transformer	17.12	37.87	21.04	39.34	150.71G	2.31M
	Mamba	14.34	32.43	18.37	35.63	51.86G	1.32M
EMGD	LSTM	14.86	32.91	19.28	36.76	196.82G	1.78M
	GRU	14.46	31.91	18.76	35.74	85.00G	1.39M
	Transformer	14.05	31.38	18.23	35.44	10.94G	0.20M
	Mamba	13.9	30.76	18.04	34.79	10.94G	0.20M

historical context and ego-motion features as conditioning information for future steps, explicitly guiding the decoder to capture the dynamic regulatory mechanism of ego-motion on pedestrian motion. This enables more accurate inference of pedestrian motion states from an ego-centric perspective. On the other hand, across all comparative experiments, models utilizing Mamba as the backbone network consistently achieve the smallest parameter counts and lowest computational costs, regardless of the decoding mechanism employed. Furthermore, models built with the proposed ego-motion-guided decoding (EMGD) mechanism generally exhibit fewer parameters and reduced computational requirements.

The effect of Mamba encoder: To further investigate the critical role of the Mamba network in motion feature extraction, we design targeted ablation experiments. In previous works, LSTM, GRU, and Transformer have been the three most prevalent sequence models for extracting pedestrian motion features and vehicle ego-motion features. In contrast, our proposed TrajMamba model adopts Mamba as the encoder to extract both types of motion features. Accordingly, we design four experimental configurations: with the Mamba network fixed as the decoder, we employ LSTM, GRU, Transformer, and Mamba as the encoders, respectively. Ablation studies are conducted on both the PIE and JAAD datasets, with the results presented in Table IV. The experimental results demonstrate that using Mamba as the encoder achieves significantly better prediction performance (i.e., lower errors) across the four metrics of ADE, FDE, ARB, and FRB compared to the other three models. This indicates that the structured state-space design of the Mamba model is more effective than recurrent neural networks and attention-based Transformer networks in extracting motion features through temporal encoding.

The effect of Mamba decoder: To further validate the critical role of Mamba as a ego-motion-guided decoder, we

TABLE IV: Ablation Results on Encoder.

Encoder	Decoder	ADE	FDE	ARB	FRB
		PIE			
LSTM	Mamba	8.13	18.23	12.28	23.58
GRU	Mamba	7.87	17.41	11.89	22.67
Transformer	Mamba	7.81	17.46	11.84	22.67
Mamba	Mamba	7.72	16.99	11.60	22.00
JAAD					
LSTM	Mamba	14.49	32.29	18.93	36.32
GRU	Mamba	14.00	31.13	18.16	35.00
Transformer	Mamba	14.38	32.37	18.88	36.89
Mamba	Mamba	13.90	30.76	18.04	34.79

TABLE V: Ablation Results on Decoder.

Encoder	Decoder	ADE	FDE	ARB	FRB
		PIE			
Mamba	LSTM	8.05	18.01	12.56	24.03
Mamba	GRU	7.98	17.90	12.33	23.59
Mamba	Transformer	7.72	17.33	11.90	22.87
Mamba	Mamba	7.72	16.99	11.60	22.00
JAAD					
Mamba	LSTM	14.54	32.44	18.95	36.40
Mamba	GRU	14.75	32.65	18.97	36.15
Mamba	Transformer	14.06	31.58	18.31	35.42
Mamba	Mamba	13.90	30.76	18.04	34.79

designed additional ablation experiments. Specifically, we fixed Mamba as the encoder and varied the decoder architecture, employing LSTM, GRU, Transformer, and Mamba as decoders, respectively. The experimental results are presented in Table V. The results demonstrate that Mamba consistently achieves the best prediction performance as the decoder on both datasets, followed by Transformer as the second-best. These results demonstrate that Mamba is effective not only as a motion encoder but also as a decoder for modeling the regulatory mechanism of ego-motion on pedestrian movement.

E. Qualitative Evaluation

In this section, we visually demonstrate the effectiveness of the proposed TrajMamba model through qualitative evaluation results. Some visualization results of prediction samples of TrajMamba on the PIE and JAAD datasets are shown in Figure 2. The first column shows the position of the pedestrian in the last observed frame, and the second to fourth columns, respectively, show the future positions at 0.5s, 1.0s, and 1.5s. The top two rows are samples from the PIE dataset, and the bottom two rows are from the JAAD dataset. These visualization results indicate that TrajMamba can effectively model the motion state of pedestrians and ego-motion and accurately predict the future trajectory of the tracked pedestrian from the egocentric perspective.



Fig. 2: Pedestrian trajectory prediction qualitative samples on the PIE and JAAD datasets. Ground truth positions in red, predictions in green. Best viewed in colour.

V. CONCLUSION

We propose a novel Mamba-based model called TrajMamba to predict the future trajectory of the tracked pedestrian from an egocentric perspective. TrajMamba innovatively employs the Mamba model to extract pedestrian motion features and the ego-motion of the ego-camera. Furthermore, an ego-motion-guided Mamba model is proposed to jointly model the relative relationship between pedestrian movement and camera movement through sequential modeling and effectively infer the future movement of pedestrians, thereby accurately predicting the future trajectory of pedestrians from the perspective of ego-motion. Experimental results on PIE and JAAD datasets validate the effectiveness and excellent performance of TrajMamba. Currently, we have not yet considered the mutual influence among pedestrians in crowds. Therefore, in future work, the prediction performance of the TrajMamba model will be improved by modeling pedestrian interaction.

REFERENCES

- [1] H. Wang, W. Zhi, G. Batista, and R. Chandra, "Pedestrian trajectory prediction using dynamics-based deep learning," in *IEEE International Conference on Robotics and Automation, ICRA 2024, Yokohama, Japan, May 13-17, 2024*. IEEE, 2024, pp. 15 068–15 075.
- [2] W. Lin, X. Zeng, C. Pang, J. Teng, and J. Liu, "Dyhgdat: Dynamic hypergraph dual attention network for multi-agent trajectory prediction," in *IEEE International Conference on Robotics and Automation, ICRA 2024, Yokohama, Japan, May 13-17, 2024*. IEEE, 2024, pp. 16 662–16 668.
- [3] A. Rasouli, M. Rohani, and J. Luo, "Bifold and semantic reasoning for pedestrian behavior prediction," in *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*. IEEE, 2021, pp. 15 580–15 590. [Online]. Available: <https://doi.org/10.1109/ICCV48922.2021.01531>
- [4] Y. Feng, T. Zhang, A. P. Sah, L. Han, and Z. Zhang, "Using Appearance to Predict Pedestrian Trajectories through Disparity-Guided Attention and Convolutional LSTM," *IEEE Transactions on Vehicular Technology*, vol. 70, no. 8, pp. 7480–7494, 2021.
- [5] Z. Zhang, Z. Ding, and R. Tian, "Decouple ego-view motions for predicting pedestrian trajectory and intention," *IEEE Trans. Image Process.*, vol. 33, pp. 4716–4727, 2024.
- [6] H. Damirchi, M. Greenspan, and A. Etemad, "Context-Aware Pedestrian Trajectory Prediction with Multimodal Transformer," in *2023 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2023, pp. 2535–2539.
- [7] X. Chen, F. Luo, F. Zhao, and Q. Ye, "Goal-Guided and Interaction-Aware State Refinement Graph Attention Network for Multi-Agent Trajectory Prediction," *IEEE Robotics and Automation Letters*, vol. 9, no. 1, pp. 57–64, 2024.
- [8] X. Zhou, X. Chen, and J. Yang, "Edge-enhanced heterogeneous graph transformer with priority-based feature aggregation for multi-agent trajectory prediction," *IEEE Trans. Intell. Transp. Syst.*, vol. 26, no. 2, pp. 2266–2281, 2025.
- [9] R. Liang, Y. Li, X. Li, Y. Tang, J. Zhou, and W. Zou, "Temporal Pyramid Network for Pedestrian Trajectory Prediction with Multi-Supervision," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, 2021, pp. 2029–2037.
- [10] D. Wang, H. Liu, N. Wang, Y. Wang, H. Wang, and S. McLoone, "SEEM: A Sequence Entropy Energy-Based Model for Pedestrian Trajectory All-Then-One Prediction," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 1, pp. 1070–1086, 2023.
- [11] X. Zhou, X. Chen, and J. Yang, "Heterogeneous hypergraph transformer network with cross-modal future interaction for multi-agent trajectory prediction," *Engineering Applications of Artificial Intelligence*, vol. 144, no. September 2024, p. 110125, 2025.

- [12] X. Chen, H. Zhang, F. Deng, J. Liang, and J. Yang, "Stochastic Non-Autoregressive Transformer-Based Multi-Modal Pedestrian Trajectory Prediction for Intelligent Vehicles," *IEEE Transactions on Intelligent Transportation Systems*, vol. 25, no. 5, pp. 3561–3574, 2024.
- [13] W. Mao, C. Xu, Q. Zhu, S. Chen, and Y. Wang, "Leapfrog Diffusion Model for Stochastic Trajectory Prediction," in *Proceedings of 2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2023, pp. 5517–5526.
- [14] W. Wu and X. Deng, "Motion Latent Diffusion for Stochastic Trajectory Prediction," in *Proceedings of 2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2024, pp. 6665–6669.
- [15] K. Guo, W. Liu, and J. Pan, "End-to-End Trajectory Distribution Prediction Based on Occupancy Grid Maps," in *Proceedings of 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, vol. 2022-June. IEEE, 2022, pp. 2232–2241.
- [16] C. Zhang, G. Zhang, Z. Zheng, and D. Lu, "Group-ptp: A pedestrian trajectory prediction method based on group features," *IEEE Trans. Multim.*, vol. 27, pp. 3527–3541, 2025.
- [17] A. Bhattacharyya, M. Fritz, and B. Schiele, "Long-term on-board prediction of people in traffic scenes under uncertainty," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4194–4202.
- [18] R. Quan, L. Zhu, Y. Wu, and Y. Yang, "Holistic LSTM for Pedestrian Trajectory Prediction," *IEEE Transactions on Image Processing*, vol. 30, pp. 3229–3239, 2021.
- [19] M. Huynh and G. Alagband, "Online Adaptive Temporal Memory with Certainty Estimation for Human Trajectory Prediction," in *Proceedings of 2023 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 2023, pp. 940–949.
- [20] S. Malla, B. Dariush, and C. Choi, "TITAN: Future Forecast Using Action Priors," in *Proceedings of 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2020, pp. 11 183–11 193.
- [21] P. Czech, M. Braun, U. Krefel, and B. Yang, "On-Board Pedestrian Trajectory Prediction Using Behavioral Features," in *Proceedings of 2022 21st IEEE International Conference on Machine Learning and Applications (ICMLA)*. IEEE, 2022, pp. 437–443.
- [22] A. Rasouli and I. Kotseruba, "PedFormer: Pedestrian Behavior Prediction via Cross-Modal Attention Modulation and Gated Multitask Learning," in *Proceedings of 2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023, pp. 9844–9851.
- [23] J. Qiu, F. P.-W. Lo, X. Gu, Y. Sun, S. Jiang, and B. Lo, "Indoor Future Person Localization from an Egocentric Wearable Camera," in *Proceedings of 2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2021, pp. 8586–8592.
- [24] A. Gu, K. Goel, and C. Ré, "Efficiently modeling long sequences with structured state spaces," in *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022.
- [25] A. Gupta, A. Gu, and J. Berant, "Diagonal state spaces are as effective as structured state spaces," in *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022*, 2022.
- [26] A. Gu and T. Dao, "Mamba: Linear-Time Sequence Modeling with Selective State Spaces," *arXiv preprint arXiv:2312.00752*, pp. 1–37, 2023.
- [27] T. Dao and A. Gu, "Transformers are ssms: Generalized models and efficient algorithms through structured state space duality," in *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net, 2024.
- [28] L. Zhu, B. Liao, Q. Zhang, X. Wang, W. Liu, and X. Wang, "Vision mamba: Efficient visual representation learning with bidirectional state space model," in *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net, 2024.
- [29] J. Park, H. Kim, K. Ko, M. Kim, and C. Kim, "Videomamba: Spatio-temporal selective state space model," in *Proceedings of 18th European Conference on Computer Vision*, vol. 15083. Springer, 2024, pp. 1–18.
- [30] H. Tang, L. Shao, Z. Zhang, L. Van Gool, and N. Sebe, "Spatial-temporal graph mamba for music-guided dance video synthesis," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. Early Access, pp. 1–11, 2025.
- [31] S. He, J. Ji, and M. Lei, "Decomposed spatio-temporal mamba for long-term traffic prediction," in *AAAI-25, Sponsored by the Association for the Advancement of Artificial Intelligence, February 25 - March 4, 2025, Philadelphia, PA, USA*, T. Walsh, J. Shah, and Z. Kolter, Eds. AAAI Press, 2025, pp. 11 772–11 780.
- [32] T. Zhang, H. Yuan, L. Qi, J. Zhang, Q. Zhou, S. Ji, S. Yan, and X. Li, "Point cloud mamba: Point cloud learning via state space model," in *AAAI-25, Sponsored by the Association for the Advancement of Artificial Intelligence, February 25 - March 4, 2025, Philadelphia, PA, USA*. AAAI Press, 2025, pp. 10 121–10 130.
- [33] J. W. Changzhi Yang, Huihui Panand, "Interact , Plan , and Go : Transformers with Social Intentions for Trajectory Prediction," *IEEE Transactions on Consumer Electronics*, vol. 71, no. 4, pp. 10 283 – 10 293, 2025.
- [34] C. Hu, R. Niu, Y. Lin, B. Yang, H. Chen, B. Zhao, and X. Zhang, "Probabilistic Trajectory Prediction of Vulnerable Road User Using Multimodal Inputs," *IEEE Transactions on Intelligent Transportation Systems*, vol. 26, no. 2, pp. 2679–2689, 2025.
- [35] H. Liu, C. Liu, F. Chang, Y. Lu, and M. Liu, "Egocentric Vulnerable Road Users Trajectory Prediction With Incomplete Observation," *IEEE Transactions on Intelligent Transportation Systems*, vol. 25, no. 10, pp. 13 694–13 705, 2024.
- [36] 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.
- [37] Y. Yao, E. M. Atkins, M. Johnson-Roberson, R. Vasudevan, and X. Du, "Bitrap: Bi-directional pedestrian trajectory prediction with multimodal goal estimation," *IEEE Robotics Autom. Lett.*, vol. 6, no. 2, pp. 1463–1470, 2021.
- [38] C. Wang, Y. Wang, M. Xu, and D. J. Crandall, "Stepwise goal-driven networks for trajectory prediction," *IEEE Robotics Autom. Lett.*, vol. 7, no. 2, pp. 2716–2723, 2022.
- [39] T. Yagi, K. Mangalam, R. Yonetani, and Y. Sato, "Future person localization in first-person videos," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7593–7602.
- [40] Y. Yao, M. Xu, C. Choi, D. J. Crandall, E. M. Atkins, and B. Dariush, "Egocentric vision-based future vehicle localization for intelligent driving assistance systems," in *International Conference on Robotics and Automation, ICRA 2019, Montreal, QC, Canada, May 20-24, 2019*. IEEE, 2019, pp. 9711–9717.
- [41] L. Neumann and A. Vedaldi, "Pedestrian and ego-vehicle trajectory prediction from monocular camera," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*. Computer Vision Foundation / IEEE, 2021, pp. 10 204–10 212.
- [42] Z. Yin, R. Liu, Z. Xiong, and Z. Yuan, "Multimodal transformer networks for pedestrian trajectory prediction," in *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI 2021, Virtual Event / Montreal, Canada, 19-27 August 2021*, Z. Zhou, Ed. ijcai.org, 2021, pp. 1259–1265. [Online]. Available: <https://doi.org/10.24963/ijcai.2021/174>
- [43] A. Rasouli, I. Kotseruba, T. Kunic, and J. K. Tsotsos, "PIE: A large-scale dataset and models for pedestrian intention estimation and trajectory prediction," in *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019*. IEEE, 2019, pp. 6261–6270.
- [44] Z. Su, G. Huang, S. Zhang, and W. Hua, "Crossmodal transformer based generative framework for pedestrian trajectory prediction," in *2022 International Conference on Robotics and Automation, ICRA 2022, Philadelphia, PA, USA, May 23-27, 2022*. IEEE, 2022, pp. 2337–2343.
- [45] A. Rasouli, "A novel benchmarking paradigm and a scale-and motion-aware model for egocentric pedestrian trajectory prediction," in *IEEE International Conference on Robotics and Automation, ICRA 2024, Yokohama, Japan, May 13-17, 2024*. IEEE, 2024, pp. 5630–5636. [Online]. Available: <https://doi.org/10.1109/ICRA57147.2024.10610614>
- [46] A. Rasouli, I. Kotseruba, and J. K. Tsotsos, "Are they going to cross? A benchmark dataset and baseline for pedestrian crosswalk behavior," in *2017 IEEE International Conference on Computer Vision Workshops, ICCV Workshops 2017, Venice, Italy, October 22-29, 2017*. IEEE Computer Society, 2017, pp. 206–213.