

# UIT-PRED: UNIVERSAL INTERMITTENT TRAJECTORY PREDICTOR FOR AUTONOMOUS DRIVING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Trajectory prediction is a fundamental component of autonomous driving, requiring models that can handle intermittent observation patterns such as variable-length histories and missing data. Existing state-of-the-art methods, however, often assume fixed-length trajectories and complete input, which challenges their applicability in real-world scenarios where sensor occlusions, communication delays, and temporal sparsity are common. Moreover, conventional approaches typically address tasks such as trajectory prediction, variable-length modeling, or missing data handling in isolation, making them less effective in multi-task settings that naturally arise in practice. To address these challenges, we propose Universal Intermittent Trajectory Predictor (UIT-Pred) that processes inputs with the time index features, which capture temporal variations to effectively adapt to diverse input patterns within the domain. Particularly, We extend recent State Space Models (SSMs) by introducing the Bidirectional Time Decay Mamba (BTD-Mamba), designed to capture dependencies both forward and backward along the sequence. By integrating a decay process, BTD-Mamba effectively analyzes trajectories while maintaining relationships under intermittent observation. Furthermore, the proposed prediction module employs state encoding to capture the underlying motion patterns in the input data and models a multimodal trajectory distribution to account for uncertainty in future predictions. These components are fused through a unified fusion module, enabling the model to jointly reason over observed dynamics and potential future behaviors. Extensive experiments on Argoverse 1 and Argoverse 2 datasets validate the effectiveness of the proposed model. By simultaneously handling prediction, variable-length observations, and missing inputs within a universal architecture, the framework proposes to meet the challenges of real-world autonomous driving systems.

## 1 INTRODUCTION

Trajectory prediction is a core challenge in autonomous driving, as safe navigation requires anticipating the future behaviors of surrounding agents under uncertain and dynamic conditions. While deep learning models Huang et al. (2025); Karim et al. (2024) has enabled significant progress, they are designed for fixed-length trajectories and complete observations. In practice, however, observations are often intermittent: variable-length sequences arise when agents enter or exit the sensor’s field of view at different times or are observed for varying durations, while missing data occurs due to sensor occlusions or communication delays. This discrepancy in the observation can degrade the performance of state-of-the-art methods unless the model explicitly handles these issues Xu & Fu (2024); Qiu et al. (2025).

Some of the approaches address the issue of variable length trajectories. Xu & Fu (2024) attribute length bias in Transformers to positional encoding and layer normalization, proposing specialized subnetworks for different sequence lengths. Li et al. (2024b) introduce a length-agnostic knowledge distillation (LaKD) module that dynamically transfers knowledge across trajectories. Qiu et al. (2025) proposes Contrastive Learning for Length Shift (CLLS), which uses contrastive learning during training to help the model learn length-invariant features and reduce the effect of varying observation lengths. Although these approaches show some effectiveness, they rely on generating multiple augmented versions of each trajectory sequence, which expands the input space and increases the overall complexity of the training.

Other methods like TransSPORTmer Capellera et al. (2024) and MS-TIP Chib et al. (2024) both apply masking techniques within transformer frameworks to address missing data, with TransSPORTmer applied to sports scenarios and MS-TIP designed for pedestrian trajectory recovery. In contrast, U2Diff Capellera et al. (2025) simultaneously reconstructs missing agent states and estimates uncertainty, focusing on the sports domain. However, these approaches depend on masking strategies to reconstruct missing observations, but this adds complexity in handling masked vs. unmasked inputs. Although effective in the sports and pedestrian domains, their applicability to autonomous driving scenarios remains largely unexplored. Additionally, State-of-the-art methods usually address prediction, variable-length observation, or missing input separately, overlooking the multi-task nature of real-world systems. Since diverse scenarios might happen in real practice, it is essential to develop a unified approach that can handle various input conditions, as illustrated in Figure. 1.

To bridge the aforementioned challenges, we propose the Universal Intermittent Trajectory Predictor (UIT-Pred), a unified architecture designed to effectively handle varying input conditions in trajectory forecasting. UIT-Pred transforms diverse input formats into a generalizable schema through time-aware input representations. Specifically, we derive two complementary temporal features from time indices: scaled timestamps to account for varying time ranges, and inter-observation intervals feature to capture the timing gaps between observations. These temporal cues enable the model to capture motion dynamics without depending on fixed time references or explicit validity masks. Building on the capabilities of recent State Space Models (SSMs), particularly the Mamba architecture, we introduce an enhanced Bidirectional Time Decay Mamba (BTD-Mamba) module, which captures sequential dependencies in both forward and backward directions across input observations. Additionally, a decay mechanism is incorporated to maintain the continuity and integrity of temporal relationships despite intermittent observations. Furthermore, in the proposed prediction module, we introduce a learnable state embedding to effectively capture the underlying dynamics of variable-length input sequences and missing observations. This embedding provides a compact yet informative representation of the agent’s motion history, maintaining temporal continuity and capturing key behavioral patterns. To further enrich this representation, we employ a cross-attention mechanism to integrate global context, including nearby agents and road topology. Finally, the enhanced state embedding is fused with the agent’s multimodal features through the proposed unified fusion module, enabling mutual learning and enhancing prediction accuracy.

Our contributions are summarized as follows: **(i)** We propose UIT-Pred, a generalizable architecture that effectively handles diverse input conditions including variable-length histories and missing input data in trajectory forecasting. **(ii)** We extend the Mamba architecture by introducing the Bidirectional Time Decay Mamba (BTD-Mamba) module, designed to extract rich spatiotemporal features from diverse forms of intermittent trajectory inputs. **(iii)** We introduce a novel prediction module that generates a learnable state embedding to capture the dynamics of observed motion patterns, which is then fused with the multimodal output to enhance trajectory prediction. **(iv)** Extensive experiments on the Argoverse 1 and Argoverse 2 benchmarks demonstrate the consistent and strong performance of our method.

## 2 RELATED WORK

### 2.1 TRAJECTORY PREDICTION

Trajectory prediction is the task of forecasting the future paths of moving agents, such as vehicles, pedestrians, or cyclist, based on a sufficiently long, fixed-length history of observed positions. In re-

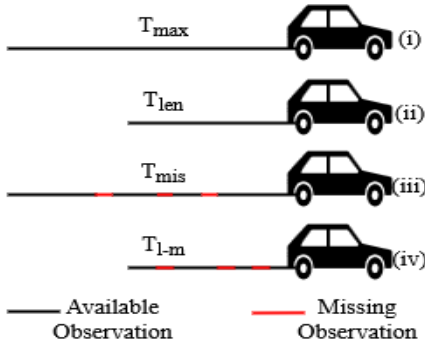


Figure 1: The following input conditions reflect scenarios commonly encountered in real-world traffic systems: (i) full observation is available; (ii) observations are available, but of variable lengths; (iii) some observation points are missing; and (iv) both variable-length and missing observations occur simultaneously.

cent years, numerous approaches have been developed to address this challenge Chen et al. (2025a), Messaoud et al. (2025). To model interactions between agents and the map, graph neural networks Wang et al. (2025a); Chen et al. (2025b) and attention-based mechanisms Xin et al. (2025); Bharilya et al. (2025); Lee et al. (2024); Huang et al. (2025) have been widely employed. Furthermore, to capture the inherent uncertainty of road agents, researchers generate multimodal predictions using GANs Wang et al. (2025b), flow-based models Liang et al. (2023), and diffusion models Capellera et al. (2025), Wang et al. (2024b), Neumeier et al. (2024). Additionally, goal-based approaches Afshar et al. (2024), Xing et al. (2025) have gained traction, where multi-modal goals are first generated through sampling or learning, followed by trajectory prediction conditioned on these goals.

Recently, the Mamba Gu & Dao (2023) framework has revived interest in state space models (SSMs) as promising alternatives to Transformers Vaswani et al. (2017), owing to their ability to reduce attention complexity and capture long-term dependencies. Mamba has shown strong potential across diverse domains, including natural language processing Zhao et al. (2025); Wang et al. (2024a) and computer vision Hatamizadeh & Kautz (2025); Yu & Wang (2025). Building on these advances, our method integrates the strengths of Mamba with Transformer architectures to achieve superior performance in the unified trajectory prediction task.

## 2.2 TRAJECTORY PREDICTION FOR LENGTH SHIFT

Trajectory prediction with variable observation lengths has received growing attention in recent years. Xu & Fu (2024) attribute length bias in Transformers to positional encoding and layer normalization, proposing specialized subnetworks for different sequence lengths. Li et al. (2024b) introduce a length-agnostic knowledge distillation (LaKD) module that dynamically transfers knowledge across trajectories. Qiu et al. (2025) proposes Contrastive Learning for Length Shift (CLLS), which uses contrastive learning during training to help the model learn length-invariant features and reduce the effect of varying observation lengths. Additionally, methods like ITPNet Li et al. (2024a), MOE Sun et al. (2022), DTO Monti et al. (2022), and SingularTrajectory Bae et al. (2024) perform instantaneous trajectory prediction by forecasting future motion based on a very short history, typically the last two time steps, but they depend on a fixed input length. In contrast, our proposed method handles variable-length observations.

## 2.3 TRAJECTORY IMPUTATION

Trajectory imputation aims to reconstruct unobserved agent states by leveraging contextual and historical motion data. Earlier work on time-series imputation has explored autoregressive RNNs for filling in missing values Cao et al. (2018). GC-VRNN Xu et al. (2023) couples a variational RNN with a spatio-temporal GNN to reconstruct missing points and forecast futures in one framework. Recently, TransPORTmer Capellera et al. (2024) applied input masking within a transformer architecture to impute missing observations, outperforming task-specific baselines in both player and ball tracking. Similarly, MS-TIP Chib et al. (2024) employed diagonal masked self-attention in transformers to recover missing data in pedestrian trajectories. U2Diff Capellera et al. (2025) introduced a unified diffusion-based model that reconstructs missing agent states while estimating state-wise uncertainty. While these methods focus on imputation, their primary applications are in sports or pedestrian settings. The challenge of handling missing data in the context of autonomous driving remains largely underexplored.

# 3 PROPOSED METHOD

**Problem Definition** In autonomous driving trajectory prediction, the goal is to forecast a target agent’s future motion based on past observations and contextual information (e.g., maps and surrounding agents). Real-world data often contain intermittent patterns, including variable-length observations and missing values. Formally, for each agent, the observed sequence is denoted as  $\tilde{X} = \{\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_{t_{\text{obs}}}\}$ , where each  $\tilde{x}_i$  contains coordinates and velocity; the length  $t_{\text{obs}}$  varies and some  $\tilde{x}_i$  may be missing. The task is to predict the future trajectory  $\hat{Y} = \{\hat{y}_1, \hat{y}_2, \dots, \hat{y}_{t_{\text{pred}}}\}$ , with each  $\hat{y}_j$  representing agent positions over a fixed horizon  $t_{\text{pred}}$ , using the available observations  $\tilde{X}$  and context.

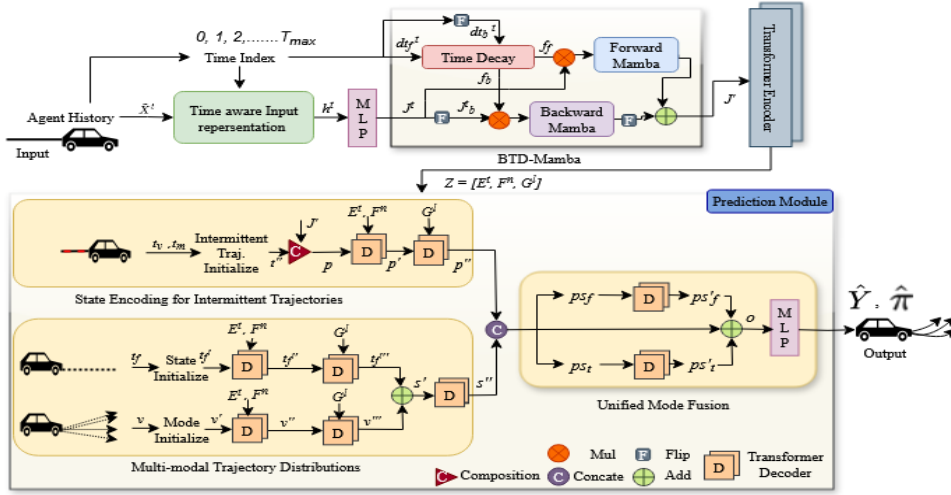


Figure 2: Illustration of our UIT-Pred framework. For simplicity, the neighbor and map encoding branches are omitted. Here,  $F^n$  represents the encoded features of neighboring agents, and  $G^l$  denotes the encoded lane information.

### 3.1 TIME-AWARE INPUT REPRESENTATION

We represent each trajectory as a sequence of temporal states with spatial and contextual features. For a given target agent, the input at each timestep  $i$  consists of concatenated state features including spatial coordinates  $t_x^{(i)}$  and velocity  $t_{\text{vel}}^{(i)}$  to form a comprehensive feature vector  $\tilde{\mathbf{X}}_i^t = [t_x^{(i)} \parallel t_{\text{vel}}^{(i)}]$  where  $\parallel$  denotes concatenation. To handle varying input lengths and missing observations, we enhance the agent’s feature state by incorporating two complementary temporal features. First, we compute a *scaled timestamps*,

$$t_i = 1 - \frac{\tau_i - \tau_{\min}}{\tau_{\max} - \tau_{\min}}, \quad t_i \in [0, 1] \quad (1)$$

where  $\tau_i$  is the original absolute timestamp (e.g.,  $\tau_i = [0, 1, \dots, t_{\text{obs}}]$ ). This maps each time index to the range  $[0, 1]$ , ensuring the model to learn temporal patterns in the observations without bias toward any specific sequence length. Second, we compute the *inter-observation interval feature*, which represents the time elapsed since the previous valid observation,

$$\Delta t_i = \tau_i - \tau_{i-1}, \quad i > 1, \quad \Delta t_1 = 0 \quad (2)$$

The feature  $\Delta t_i$  allow the model to reason about missing observations and We explicitly tell the model how much time has passed since the previous observation., without explicitly relying on binary validity masks. Thus, the final input at each timestep is represented as,

$$\mathbf{h}_i^t = [t_x^{(i)} \parallel t_{\text{vel}}^{(i)} \parallel t_{\text{norm}}^{(i)} \parallel \Delta t^{(i)}]_{i=1}^{t_{\text{obs}}} \quad (3)$$

This formulation integrates spatial-temporal states, scaled timestamps to handle varying time ranges, and inter-observation interval feature to represent the timing differences in the observation.

For each neighboring agent, the input representation is constructed similarly to that of the target agent, using all available timestamps within a fixed observation window to form the feature vector  $\mathbf{h}_i^n$ . For lane segment points, the input  $\mathbf{h}^l$  combines geometric and visibility information, following the design in Zhang et al. (2024). (refer to Appendix A.1 for further details).

### 3.2 BIDIRECTIONAL TIME-DECAY MAMBA (BTD-MAMBA)

Mamba blocks inherently support variable-length sequences through recurrent state-space updates. Building on this, we introduce BTM-Mamba, as shown in Figure. 2, an enhanced state-space model

designed to handle both variable-length and missing trajectory data. It extends Mamba by incorporating a time decay mechanism that modulates hidden states based on inter-observation intervals. Initially, we compute forward and backward inter-arrival times as,

$$\mathbf{dt}_f^t = [\Delta t_0, \Delta t_1, \dots, \Delta t_{T_{\text{obs}}-1}], \quad \mathbf{dt}_b^t = \text{flip}(\mathbf{dt}_f^t), \quad \mathbf{dt}^{(f,b)} = [\mathbf{dt}_f^t, \mathbf{dt}_b^t] \quad (4)$$

The concatenated inter-arrival times  $\mathbf{dt}^{(f,b)}$  encode both forward and backward time gaps between observations. We then project this into a scaling feature space using,

$$[f_t, f_b] = \frac{1}{\exp\left(\phi_s\left(\mathbf{dt}^{(f,b)}; W_s\right)\right)} \quad (5)$$

Here,  $\phi_s(\cdot)$  is a projection function parameterized by weights  $W_s$ , implemented as a multi-layer perceptron (MLP) with ReLU activation. Equation 4 is designed to calculate the distance from the last observation to the current time step, which helps quantify the influence of temporal gaps, particularly when dealing with complex missing patterns. The key insight is that the influence of a variable that has been missing for a period decreases over time. Therefore, in Equation 5 we utilize a negative exponential function combined with ReLU to ensure that the influence decays monotonically within a reasonable range between 0 and 1. Moreover, we apply Mamba bidirectionally with a time-decay mechanism, effectively capturing irregular temporal intervals and modeling temporal gaps by leveraging inter-arrival times in both directions. The embedded sequence  $\mathbf{J}^t = \{\mathbf{j}_1, \mathbf{j}_2, \dots, \mathbf{j}_{t_{\text{obs}}}\}$ , where each  $\mathbf{j}_i$  is generated by passing  $\mathbf{h}_i^t$  through an MLP, is processed both in its original order and reversed order  $\mathbf{J}_b^t = \{\mathbf{j}_{t_{\text{obs}}}, \dots, \mathbf{j}_2, \mathbf{j}_1\}$  by the revised Mamba block as described below,

$$\mathbf{J}' = (\mathbf{J}^t \odot \mathbf{f}_f) * C_{\text{forw}} + \text{Flip}\left((\mathbf{J}_b^t \odot \mathbf{f}_b) * C_{\text{back}}\right) \quad (6)$$

where  $\odot$  is element-wise multiplication,  $*$  is convolution,  $\text{Flip}(\cdot)$  reverses the sequence, and  $C_{\text{forw}}, C_{\text{back}}$  are learnable convolution kernels for forward and backward directions, respectively.

**Interaction Representation** The neighboring agent features  $\mathbf{h}_i^n$  are first embedded using an MLP to produce  $\mathbf{J}^n$ . These embeddings are then passed through Mamba blocks and subsequently refined via residual layers with skip connections and normalization, as expressed below,

$$\mathbf{J}^n = \text{MLP}(\mathbf{h}_i^n), \quad \mathbf{F} = \text{MambaBlocks}(\mathbf{J}^n), \quad \mathbf{F}' = \text{Norm}(\mathbf{F} + \mathbf{J}^n) \quad (7)$$

To encode lane segment features  $\mathbf{h}^l$ , we adopt a PointNet-based polyline encoder Zhang et al. (2024), producing map embeddings  $\mathbf{G}' = \text{PNEncoder}(\mathbf{h}^l)$ . To capture interactions among the target agent, neighboring agents, and the map, their embeddings are concatenated and processed by a Transformer encoder,

$$\mathbf{I} = [\mathbf{J}' \parallel \mathbf{F}' \parallel \mathbf{G}'], \quad \mathbf{Z} = \text{TransformerEncoder}(\mathbf{I}) \quad (8)$$

where  $\mathbf{Z} = [\mathbf{E}^t, \mathbf{F}^n, \mathbf{G}^l]$  captures joint contextual representations for predicting the next trajectory.

### 3.3 PREDICTION MODULE

**State Encoding for Intermittent Trajectories** In our framework, to effectively capture the state and temporal dynamics of trajectories amid variable sequence lengths and missing observations, we introduce a *variable-step temporal representation* that accounts for unobserved, varying sequence lengths, alongside a *missing-step temporal embedding* to handle missing data. Given a sequence length  $t_{\text{len}}$ , we construct a time vector starting from a fixed maximum time  $T_{\text{max}}$  and decreasing to  $T_{\text{max}} - t_{\text{len}} + 1$ . The vector  $t_v$  represents the variable-length unobserved time steps, while  $t_m$  denotes the count of missing time steps  $m_i$ ,

$$\mathbf{t}_v = [T_{\text{max}}, \dots, T_{\text{max}} - t_{\text{len}} + 1], \quad \mathbf{t}_m = [m_1, \dots, m_{t_m}], \quad \mathbf{t} = [t_v, t_m] \quad (9)$$

These time values  $\mathbf{t}$  are normalized and scaled before being embedded via an MLP to produce learnable temporal features,

$$\mathbf{t}_{\text{scaled}} = 0.1 \times \mathbf{t} + 0.1, \quad \mathbf{t}' = \text{MLP}(\mathbf{t}_{\text{scaled}}) \quad (10)$$

The resulting embeddings  $\mathbf{t}'$  are added to repeated latent features of the target agent, enriched with state information via a GRU,

$$\mathbf{t}'' = \mathbf{t}' \oplus \text{GRU}(\mathbf{E}^t), \quad (11)$$

where  $\oplus$  denotes broadcasting addition across time and  $\mathbf{t}'' = [v'_t, m'_t]$  represents the final variable-step temporal embedding  $v'_t$  and missing-step temporal embedding  $m'_t$ .

Moreover, we integrate the observed sequence features with the variable-step temporal representation and missing-step temporal embedding to create a unified observed sequence, enabling seamless learning across all components. Initially, a ternary mask  $\mathbf{M} \in \{0, 1, 2\}^{t_{\max}}$  is constructed,

$$\mathbf{M}_{i,t} = \begin{cases} 0 & \text{if time step } t \text{ is observed,} \\ 1 & \text{if time step } t \in \mathbf{v}'_t, \\ 2 & \text{if time step } t \in \mathbf{m}'_t. \end{cases} \quad (12)$$

Using this mask, we construct the tensor  $\mathbf{p}$ , which contains observed data, final variable-step temporal embedding and missing-step temporal embedding,

$$\mathbf{p}[i,t] = \begin{cases} \mathbf{J}'[i,j], & \text{if } \mathbf{M}[i,t] = 0 \\ \mathbf{t}'_m[i,k], & \text{if } \mathbf{M}[i,t] = 1 \\ \mathbf{t}'_{\text{imp}}[i,l], & \text{if } \mathbf{M}[i,t] = 2 \end{cases} \quad (13)$$

Finally, the reconstructed full past-length sequence  $\mathbf{p}$  is generated, enriched with agent information  $[E^t, F^n]$  and lane information  $[G^l]$  through cross-attention,

$$\mathbf{p}' = \mathcal{D}(\mathbf{p}, [E^t, F^n]), \quad \mathbf{p}'' = \mathcal{D}(\mathbf{p}', [G^l]) \quad (14)$$

where  $\mathcal{D}$  is the Transformer decoder for cross-attention, and  $\mathbf{p}''$  encodes agent and lane information.

**multi-modal trajectory distributions** To generate a multimodal future trajectory for an agent, we require both mode and state information. The initial mode vector  $v$  is embedded and combined with the target agent’s encoding to form  $v'$ , which is then refined via cross-attention with agent  $[E^t, F^n]$  and lane  $[G^l]$  information,

$$\mathbf{v}' = \mathbf{v} + \mathbf{E}^t, \quad \mathbf{v}'' = \mathcal{D}(\mathbf{v}', [E^t, F^n]), \quad \mathbf{v}''' = \mathcal{D}(\mathbf{v}'', [G^l]) \quad (15)$$

where  $\mathbf{v}'''$  represents the final mode vector. To generate the state vector  $\mathbf{t}_f$ , a normalized time embedding is constructed and fused with the GRU hidden states,

$$\mathbf{t}'_f = \text{MLP}(0.1 \cdot t + 0.1) \oplus \text{GRU}[E^t], \quad t = 1, \dots, t_{\text{pred}} \quad (16)$$

The  $\mathbf{t}'_f$ -enriched states are refined via cross-attention with agent  $[E^t, F^n]$  and lane  $[G^l]$  information,

$$\mathbf{t}''_f = \mathcal{D}(\mathbf{t}'_f, [E^t, F^n]), \quad \mathbf{t}'''_f = \mathcal{D}(\mathbf{t}''_f, [G^l]) \quad (17)$$

The final  $\mathbf{t}'''_f$  integrates both agent and lane information. Furthermore, the embedding of multimodal future trajectory  $\mathbf{s}'$  is generated by combining the mode vector  $\mathbf{v}'''$  and state vectors  $\mathbf{t}'''_f$  and refined with  $Z = [E^t, F^n, G^l]$ ,

$$\mathbf{s}' = \mathbf{v}''' \oplus \mathbf{t}'''_f \quad \mathbf{s}'' = \mathcal{D}(\mathbf{s}', Z) \quad (18)$$

where  $\mathbf{s}''$  is the refined future trajectory embedding after processing.

**Unified Mode Fusion for Prediction** To capture relationships between past and future behaviors, we fuse representations  $p''$  and  $s''$  into  $ps_f$  by concatenating along the feature dimension  $\|_{\text{feat}}$  (preserving temporal resolution) and along the temporal dimension  $\|_{\text{time}}$  (aligning sequences),

$$ps_f = [p'' \|_{\text{feat}} s''], \quad ps_t = [p'' \|_{\text{time}} s''] \quad (19)$$

These fused mode representations are combined to form comprehensive tensors  $ps_f$  and  $ps_t$ , which are subsequently processed using two cross-attention mechanisms,

$$ps'_f = \mathcal{D}(ps_f, Z), \quad ps'_t = \mathcal{D}(ps_t, Z) \quad (20)$$

Finally, the outputs from the different blocks are summarized,

$$\mathbf{o} = ps_t \oplus \text{reshape}_{\text{time}}(ps'_f) \oplus ps'_t \quad (21)$$

The tensor  $\mathbf{o}[-t_{\text{pred}} : ]$  is used for downstream multimodal trajectory prediction, while  $\mathbf{o}[t_v, t_m]$  is used to predict the corresponding observations,

$$\hat{Y}, \hat{\pi} = \text{MLP}(\text{output}[-t_{\text{pred}} : ]), \quad \hat{x}_v, \hat{x}_m = \text{MLP}(\text{output}[t_v, t_m]), \quad (22)$$

where  $\hat{Y}$  denotes the predicted future trajectories,  $\hat{\pi}$  represents the associated mode probabilities, and  $\hat{x}_v, \hat{x}_m$  are the reconstructed variable-length steps and missing observations, respectively.

Table 1: Performance under different observation scenarios for Argoverse 2 dataset.

Model	I/N Scenarios	minADE <sub>1</sub>	minFDE <sub>1</sub>	MR <sub>1</sub>	minADE <sub>6</sub>	minFDE <sub>6</sub>	MR <sub>6</sub>
DeMo-Orig	Missing + Var.	2.1672	4.9458	0.6457	0.8889	1.5834	0.2065
	Var. Obs	2.0281	4.5055	0.6329	0.8177	1.5054	0.1971
	Missing Only	1.6637	4.1427	0.5944	0.6837	1.3234	0.1666
	Full Obs	1.6041	4.0521	0.5889	0.6625	1.2999	0.1623
DeMo-RSD	Missing + Var.	1.6732	4.1777	0.5960	0.6812	1.33277	0.1683
	Var. Obs	1.6613	4.1589	0.5955	0.6769	1.3261	0.1663
	Missing Only	1.6747	4.1703	0.5956	0.6863	1.3419	0.1697
	Full Obs	1.6341	4.1049	0.5914	0.6719	1.3187	0.1646
Forecast-mae-Orig	Missing + Var.	2.1137	5.0919	0.6735	0.8291	1.5985	0.2130
	Var. Obs	2.0186	4.8589	0.6406	0.7892	1.5182	0.2034
	Missing Only	2.2330	5.1860	0.6649	0.8524	1.6087	0.2195
	Full Obs	1.8165	4.5536	0.6218	0.7244	1.4273	0.1877
Forecast-mae-RSD	Missing + Var.	1.8246	4.5432	0.6201	0.7335	1.4406	0.1908
	Var. Obs	1.8146	4.5292	0.6204	0.7292	1.4354	0.1902
	Missing Only	1.8302	4.5558	0.6224	0.7346	1.4423	0.1918
	Full Obs	1.8098	4.5198	0.6200	0.7279	1.4332	0.1905
Our	Missing + Var.	1.5882	3.9402	0.5717	0.6562	1.2326	0.1551
	Var. Obs	1.5749	3.9281	0.5693	0.6490	1.2316	0.1525
	Missing Only	1.5711	3.9059	0.5692	0.6491	1.2481	0.1527
	Full Obs	<b>1.5508</b>	<b>3.8784</b>	<b>0.5677</b>	<b>0.6436</b>	<b>1.2455</b>	<b>0.1526</b>

**Model Training.** To supervise the predicted trajectory and its confidence, we employ the Huber loss for trajectory regression, denoted as  $L_{\text{reg}}$ , along with the cross-entropy loss for confidence classification, denoted as  $L_{\text{cls}}$ . Additionally, the embeddings for variable-step temporal embedding and missing-step temporal embedding are supervised using the  $L_{\text{u-rg}}$  and  $L_{\text{m-rg}}$  losses, respectively, both of which employ the same regression criterion. Furthermore, an endpoint loss  $L_{\text{et}}$  is incorporated for all agents, applying the same regression loss function to improve endpoint accuracy. The model is trained end-to-end by combining all the losses,

$$L_{\text{total}} = L_{\text{reg}} + L_{\text{cls}} + L_{\text{u-rg}} + L_{\text{m-rg}} + L_{\text{et}} \quad (23)$$

All loss terms are weighted equally, and  $L_{\text{total}}$  represents the overall training loss.

## 4 EXPERIMENTS

### 4.1 EXPERIMENTAL SETTINGS

**Datasets & Evaluation Metric** Our method is evaluated on two benchmark datasets such as Argoverse 1 Chang et al. (2019) with 323,557 sequences including 2 seconds of past data and 3 seconds of future prediction, and Argoverse 2 Wilson et al. (2021) with 250,000 scenes providing 5 seconds of observation and 6 seconds of prediction. We assess our approach using the standard metrics  $MinADE_k$ ,  $MinFDE_k$ , and  $MR_k$ , with  $k$  values of 1 and 6 commonly adopted as benchmarks.

**Implementation Details** Detailed training settings are included in the Appendix A.3 section.

### 4.2 RESULT AND ANALYSIS

**Performance under Different Observation Settings.** We evaluate the overall performance of the proposed method under various input conditions on the Argoverse 2 dataset, as presented in Table 1. The label *Missing + Var* indicates scenarios where inputs have variable lengths and contain missing data. *Var. Obs* refers to inputs with variable lengths only, without missing data. *Missing Only* denotes inputs that contain missing data. Finally, *Full Obs* represents complete inputs with fixed

Table 2: Comparison of different methods on Argoverse 2 and Argoverse 1 with variable-length observations. The best results are highlighted in **bold**.

Dataset	Method	minADE <sub>1</sub>	minFDE <sub>1</sub>	MR <sub>1</sub>	minADE <sub>6</sub>	minFDE <sub>6</sub>	MR <sub>6</sub>
Argoverse 2	HiVT-Orig	2.5502	6.5586	0.7455	1.0561	2.1093	0.3275
	HiVT-RM	2.2848	6.0548	0.7249	0.9457	1.9283	0.2994
	HiVT-DTO	2.2769	6.0548	0.7275	0.9324	1.8946	0.2903
	HiVT-FLN	2.2786	6.0464	0.7240	0.9287	1.8838	0.2891
	HiVT-LaKD	2.2066	5.8769	0.7161	0.9183	1.8686	0.2791
	QCNet-Orig	2.1006	5.2219	0.6299	0.8339	1.3849	0.1884
	QCNet-RM	1.7452	4.4404	0.5957	0.7508	1.3184	0.1671
	QCNet-DTO	1.7713	4.4900	0.5979	0.7454	1.2924	0.1671
	QCNet-FLN	1.6940	4.2373	0.5808	0.7370	1.2595	0.1596
	QCNet-LaKD	1.6574	4.1505	0.5753	0.7258	1.2420	0.1555
	<b>our</b>	<b>1.5749</b>	<b>3.9281</b>	<b>0.5693</b>	<b>0.6490</b>	<b>1.2316</b>	<b>0.1525</b>
Argoverse 1	HiVT-Orig	1.4733	3.1834	0.5267	0.7255	1.0740	0.1124
	HiVT-RM	1.4189	3.0599	0.5104	0.7070	1.0447	0.1053
	HiVT-DTO	1.3999	3.0262	0.5056	0.7032	1.0350	0.1039
	HiVT-FLN	1.4011	3.0288	0.5051	0.7026	1.0325	0.1033
	HiVT-LaKD	1.3317	2.8799	0.4901	0.6807	0.9864	0.0928
	QCNet-Orig	1.1656	2.4021	0.3860	0.5791	0.7399	0.0734
	QCNet-RM	1.0995	2.2550	0.3630	0.5684	0.7115	0.0703
	QCNet-DTO	1.0708	2.2303	0.3563	0.5418	0.6848	0.0671
	QCNet-FLN	1.0631	2.2083	0.3579	0.5411	0.6680	0.0671
	QCNet-LaKD	0.9982	2.0718	0.3439	0.5240	0.6581	0.0640
	<b>our</b>	<b>0.8086</b>	<b>1.7343</b>	<b>0.2871</b>	<b>0.3693</b>	<b>0.5901</b>	<b>0.0510</b>

length. The models *DeMo-Orig* and *Forecast-mae-Orig* are trained exclusively on fully observed, fixed-length inputs and are evaluated across all input conditions to examine their generalization capability. In contrast, *DeMo-RSD* and *Forecast-mae-RSD* are trained on inputs with randomly assigned sequence lengths, where certain time steps are also randomly dropped during training. This training setup is referred to as *RSD* (Random Sequence Drop).

The performance of *DeMo-Orig* and *Forecast-mae-Orig* degrades significantly under the *Missing + Var*, *Var: Obs*, and *Missing Only* settings, demonstrating their limitations in handling missing or variable-length inputs. In contrast, *DeMo-RSD* and *Forecast-mae-RSD*, although trained with the *RSD* strategy, do not show notable improvements in these challenging scenarios. This suggests that *RSD* alone is not sufficient to ensure adaptability across diverse input conditions. The proposed model, *UIT-Pred*, demonstrates promising results across all input conditions. (See Appendix A.2 for comparisons on the Argoverse 1 dataset.)

**Performance under Different Observations Lengths.** The performance of the proposed approach is evaluated in Table 2 with different observation lengths. The *HiVT-Orig* and *QCNet-Orig* models refer to the original versions of HiVT and QCNet trained using fixed-length observed trajectories as input. In contrast, *HiVT-RM* and *QCNet-RM* introduce random masking to the observed trajectories during training to simulate inputs of varying lengths. Variants such as *HiVT-DTO*, *HiVT-FLN*, and *HiVT-LaKD* represent configurations that use HiVT as the backbone, combined with the DTO, FlexiLength Network (FLN), and Length-agnostic Knowledge Distillation (LaKD) modules, respectively. Similarly, *QCNet-DTO*, *QCNet-FLN*, and *QCNet-LaKD* use QCNet as the backbone in combination with the same modules.

The results demonstrate that the proposed model consistently outperforms all baseline methods, including *QCNet-FLN*, *QCNet-LaKD*, *HiVT-FLN*, and *HiVT-LaKD*, across both data sets. As these baselines are specifically designed to handle variable-length inputs, this comparison highlights the generalizability of our approach. Furthermore, our method also surpasses *HiVT-Orig* and *QCNet-*



Table 3: Component Study of Proposed Model for Argoverse 2 dataset

ID	TAIR	BTD-M.	PM	minADE <sub>1</sub>	minFDE <sub>1</sub>	MR <sub>1</sub>	minADE <sub>6</sub>	minFDE <sub>6</sub>	MR <sub>6</sub>
1				1.8528	4.6447	0.6442	0.8189	1.4668	0.1915
2	✓			1.8071	4.3175	0.6257	0.7538	1.4381	0.1830
3		✓		1.7847	4.2176	0.6137	0.7381	1.4169	0.1798
4			✓	1.7516	4.0163	0.6149	0.7257	1.4037	0.1705
5	✓	✓		1.6928	3.9826	0.6037	0.7081	1.3569	0.1629
6	✓		✓	1.6683	3.9471	0.5973	0.6822	1.3244	0.1648
7		✓	✓	1.6528	3.9714	0.5901	0.6962	1.3062	0.1631
8	✓	✓	✓	1.5882	3.9402	0.5717	0.6562	1.2326	0.1551

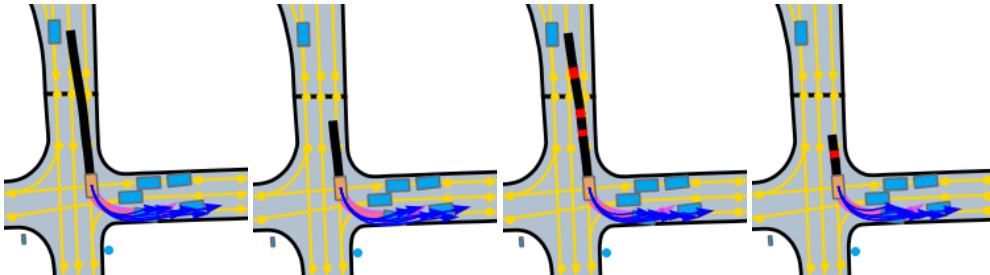


Figure 3: Qualitative results of the proposed model with varying input observations. Black: observed history; red: missing points; magenta: ground truth future; blue: predicted future trajectory.

*Orig*, reinforcing the importance of trajectory prediction frameworks tailored for variable-length observations. Despite the use of random masking in *HiVT-RM* and *QCNNet-RM*, our model still achieves superior performance, demonstrating the advantages of our structured design under varying input lengths. Overall, proposed method achieves state-of-the-art results across multiple configurations, confirming its effectiveness and adaptability.

**Ablation of Each Component** Table 3 presents a component study of the proposed model, evaluating the contributions of the Time-Aware Input Representation (TAIR), BTD-Mamba (BTD-M.), and the Predictor Module (PM). When BTD-Mamba is not used, the model depended only on the forward Mamba module without Time Decay (TD) and when the PM is excluded, separate MLPs are used for each prediction. The model performs the worst when none of the components are used (ID-1), highlighting their necessity. Introducing at least one component (ID-2, 3, 4) leads to noticeable performance gains, as each individual module provides valuable information. When any two components are combined (ID-5, 6, 7), the model benefits from their complementary strengths, further improving performance. Finally, the best results are achieved when all three components are used together (ID-8), highlighting their benefit of integration.

**Qualitative Result** Figure 3 shows the proposed model’s performance during a left-turn maneuver at an intersection in the Argoverse 2 dataset. The model effectively handles all input conditions, aligning well with the real-world requirements of autonomous driving. Additional qualitative results across diverse driving scenarios are provided in Appendix A.3.

## 5 CONCLUSION

In this work, we introduce UIT-Pred, a universal architecture that addresses real-world challenges in autonomous driving by handling diverse input types for trajectory prediction. We propose a time-aware input representation that helps the model focus on motion dynamics across diverse input conditions rather than absolute durations. Furthermore, we extend a state space model to develop the BTD-Mamba module and introduce a novel predictor, jointly capturing complex temporal dynamics to enhance trajectory prediction accuracy. Comprehensive experiments on the Argoverse 1 and Argoverse 2 datasets demonstrate effectiveness of our approach.

## REFERENCES

- 486  
487  
488 Sepideh Afshar, Nachiket Deo, Akshay Bhagat, Titas Chakraborty, Yunming Shao, Balarama Raju  
489 Buddharaju, Adwait Deshpande, and Henggang Cui. Pbp: Path-based trajectory predic-  
490 tion for autonomous driving. In *2024 IEEE International Conference on Robotics and Automation*  
491 *(ICRA)*, pp. 12927–12934. IEEE, 2024.
- 492  
493 Inhwan Bae, Young-Jae Park, and Hae-Gon Jeon. Singulartrajectory: Universal trajectory predictor  
494 using diffusion model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
*Pattern Recognition*, pp. 17890–17901, 2024.
- 495  
496 Vibha Bharilya, Ashok Arora, and Neetesh Kumar. Self-supervised transformer for trajectory pre-  
497 diction using noise imputed past trajectory. *IEEE Transactions on Intelligent Transportation*  
498 *Systems*, 2025.
- 499  
500 Wei Cao, Dong Wang, Jian Li, Hao Zhou, Lei Li, and Yitan Li. Brits: Bidirectional recurrent  
imputation for time series. *Advances in neural information processing systems*, 31, 2018.
- 501  
502 Guillem Capellera, Luis Ferraz, Antonio Rubio, Antonio Agudo, and Francesc Moreno-Noguer.  
Transportmer: A holistic approach to trajectory understanding in multi-agent sports. In *Proceed-*  
503 *ings of the asian conference on computer vision*, pp. 1652–1670, 2024.
- 504  
505 Guillem Capellera, Antonio Rubio, Luis Ferraz, and Antonio Agudo. Unified uncertainty-aware  
506 diffusion for multi-agent trajectory modeling. In *Proceedings of the Computer Vision and Pattern*  
507 *Recognition Conference*, pp. 22476–22486, 2025.
- 508  
509 Ming-Fang Chang, John Lambert, Patsorn Sangkloy, Jagjeet Singh, Slawomir Bak, Andrew Hart-  
510 nett, De Wang, Peter Carr, Simon Lucey, Deva Ramanan, et al. Argoverse: 3d tracking and  
511 forecasting with rich maps. In *Proceedings of the IEEE/CVF conference on computer vision and*  
*pattern recognition*, pp. 8748–8757, 2019.
- 512  
513 Kai Chen, Xiaodong Zhao, Yujie Huang, Guoyu Fang, Xiao Song, Ruiping Wang, and Ziyuan Wang.  
Socialmoif: Multi-order intention fusion for pedestrian trajectory prediction. In *Proceedings of*  
514 *the Computer Vision and Pattern Recognition Conference*, pp. 22465–22475, 2025a.
- 515  
516 Wangxing Chen, Haifeng Sang, Jinyu Wang, and Zishan Zhao. Dstignc: Deformable spatial-  
517 temporal interaction graph convolution network for pedestrian trajectory prediction. *IEEE Trans-*  
518 *actions on Intelligent Transportation Systems*, 2025b.
- 519  
520 Jie Cheng, Xiaodong Mei, and Ming Liu. Forecast-mae: Self-supervised pre-training for motion  
521 forecasting with masked autoencoders. In *Proceedings of the IEEE/CVF International Conference*  
*on Computer Vision*, pp. 8679–8689, 2023.
- 522  
523 Pranav Singh Chib, Achintya Nath, Paritosh Kabra, Ishu Gupta, and Pravendra Singh. Ms-tip: impu-  
524 tation aware pedestrian trajectory prediction. In *International Conference on Machine Learning*,  
pp. 8389–8402. PMLR, 2024.
- 525  
526 Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv*  
527 *preprint arXiv:2312.00752*, 2023.
- 528  
529 Ali Hatamizadeh and Jan Kautz. Mambavision: A hybrid mamba-transformer vision backbone.  
530 In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 25261–25270,  
2025.
- 531  
532 Yizhou Huang, Yihua Cheng, and Kezhi Wang. Trajectory mamba: Efficient attention-mamba fore-  
533 casting model based on selective ssm. In *Proceedings of the Computer Vision and Pattern Recog-*  
*niton Conference*, pp. 12058–12067, 2025.
- 534  
535 Rezaul Karim, Soheil Mohamad Alizadeh Shabestary, and Amir Rasouli. Destine: Dynamic goal  
536 queries with temporal transductive alignment for trajectory prediction. In *2024 IEEE Interna-*  
537 *tional Conference on Robotics and Automation (ICRA)*, pp. 2230–2237. IEEE, 2024.
- 538  
539 Seongju Lee, Junseok Lee, Yeonguk Yu, Taeri Kim, and Kyoobin Lee. Mart: Multiscale relational  
transformer networks for multi-agent trajectory prediction. In *European Conference on Computer*  
*Vision*, pp. 89–107. Springer, 2024.

- 540 Rongqing Li, Changsheng Li, Yuhang Li, Hanjie Li, Yi Chen, Ye Yuan, and Guoren Wang. Itpnet:  
541 Towards instantaneous trajectory prediction for autonomous driving. In *Proceedings of the 30th*  
542 *ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 1643–1654, 2024a.
- 543  
544 Yuhang Li, Changsheng Li, Ruilin Lv, Rongqing Li, Ye Yuan, and Guoren Wang. Lakd: Length-  
545 agnostic knowledge distillation for trajectory prediction with any length observations. *Advances*  
546 *in Neural Information Processing Systems*, 37:28720–28744, 2024b.
- 547 Rongqin Liang, Yuanman Li, Jiantao Zhou, and Xia Li. Stglow: A flow-based generative framework  
548 with dual-graphormer for pedestrian trajectory prediction. *IEEE transactions on neural networks*  
549 *and learning systems*, 2023.
- 550  
551 Kaouther Messaoud, Matthieu Cord, and Alexandre Alahi. Towards generalizable trajectory pre-  
552 diction using dual-level representation learning and adaptive prompting. In *Proceedings of the*  
553 *Computer Vision and Pattern Recognition Conference*, pp. 27564–27574, 2025.
- 554  
555 Alessio Monti, Angelo Porrello, Simone Calderara, Pasquale Coscia, Lamberto Ballan, and Rita  
556 Cucchiara. How many observations are enough? knowledge distillation for trajectory forecasting.  
557 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.  
558 6553–6562, 2022.
- 559  
560 Marion Neumeier, Sebastian Dorn, Michael Botsch, and Wolfgang Utschick. Reliable trajectory  
561 prediction and uncertainty quantification with conditioned diffusion models. In *Proceedings of*  
562 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3461–3470, 2024.
- 563  
564 Ruiqi Qiu, Jun Gong, Xinyu Zhang, Siqi Luo, Bowen Zhang, and Yi Cen. Adapting to observation  
565 length of trajectory prediction via contrastive learning. In *Proceedings of the Computer Vision*  
566 *and Pattern Recognition Conference*, pp. 1645–1654, 2025.
- 567  
568 Jianhua Sun, Yuxuan Li, Liang Chai, Hao-Shu Fang, Yong-Lu Li, and Cewu Lu. Human trajectory  
569 prediction with momentary observation. In *Proceedings of the IEEE/CVF conference on computer*  
570 *vision and pattern recognition*, pp. 6467–6476, 2022.
- 571  
572 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,  
573 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2017.
- 574  
575 Chengyue Wang, Haicheng Liao, Bonan Wang, Yanchen Guan, Bin Rao, Ziyuan Pu, Zhiyong Cui,  
576 Cheng-Zhong Xu, and Zhenning Li. Nest: A neuromodulated small-world hypergraph trajectory  
577 prediction model for autonomous driving. In *Proceedings of the AAAI Conference on Artificial*  
578 *Intelligence*, volume 39, pp. 808–816, 2025a.
- 579  
580 Jingyuan Wang, Yujing Lin, and Yudong Li. Gtg: Generalizable trajectory generation model for  
581 urban mobility. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp.  
582 834–842, 2025b.
- 583  
584 Junxiong Wang, Daniele Paliotta, Avner May, Alexander Rush, and Tri Dao. The mamba in the  
585 llama: Distilling and accelerating hybrid models. *Advances in Neural Information Processing*  
586 *Systems*, 37:62432–62457, 2024a.
- 587  
588 Yixiao Wang, Chen Tang, Lingfeng Sun, Simone Rossi, Yichen Xie, Chensheng Peng, Thomas  
589 Hannagan, Stefano Sabatini, Nicola Poerio, Masayoshi Tomizuka, et al. Optimizing diffusion  
590 models for joint trajectory prediction and controllable generation. In *European Conference on*  
591 *Computer Vision*, pp. 324–341. Springer, 2024b.
- 592  
593 Benjamin Wilson, William Qi, Tanmay Agarwal, John Lambert, Jagjeet Singh, Siddhesh Khandel-  
wal, Bowen Pan, Ratnesh Kumar, Andrew Hartnett, Jhony Kaesemodel Pontes, et al. Argoverse  
2: Next generation datasets for self-driving perception and forecasting. In *Thirty-fifth Conference*  
*on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021.
- Guipeng Xin, Duanfeng Chu, Liping Lu, Zejian Deng, Yuang Lu, and Xigang Wu. Multi-agent  
trajectory prediction with difficulty-guided feature enhancement network. *IEEE Robotics and*  
*Automation Letters*, 2025.

594 Zebin Xing, Xingyu Zhang, Yang Hu, Bo Jiang, Tong He, Qian Zhang, Xiaoxiao Long, and Wei  
595 Yin. Goalflow: Goal-driven flow matching for multimodal trajectories generation in end-to-end  
596 autonomous driving. In *Proceedings of the Computer Vision and Pattern Recognition Conference*,  
597 pp. 1602–1611, 2025.

598  
599 Yi Xu and Yun Fu. Adapting to length shift: Flexilength network for trajectory prediction. In *Pro-  
600 ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15226–  
601 15237, 2024.

602 Yi Xu, Armin Bazarjani, Hyung-gun Chi, Chiho Choi, and Yun Fu. Uncovering the missing pattern:  
603 Unified framework towards trajectory imputation and prediction. In *Proceedings of the IEEE/CVF  
604 Conference on Computer Vision and Pattern Recognition*, pp. 9632–9643, 2023.

605  
606 Weihao Yu and Xinchao Wang. Mambaout: Do we really need mamba for vision? In *Proceedings  
607 of the Computer Vision and Pattern Recognition Conference*, pp. 4484–4496, 2025.

608 Bozhou Zhang, Nan Song, and Li Zhang. Demo: Decoupling motion forecasting into directional  
609 intentions and dynamic states. In *The Thirty-eighth Annual Conference on Neural Information  
610 Processing Systems*, 2024. URL <https://openreview.net/forum?id=rbtnRsixSN>.

611  
612 Han Zhao, Min Zhang, Wei Zhao, Pengxiang Ding, Siteng Huang, and Donglin Wang. Cobra:  
613 Extending mamba to multi-modal large language model for efficient inference. In *Proceedings of  
614 the AAAI Conference on Artificial Intelligence*, volume 39, pp. 10421–10429, 2025.

615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647

## 648 A APPENDIX

### 649 A.1 EXPERIMENTAL SETTINGS

650 **Input Representation for Neighbours and Lane information.** Using the vectorized representa-  
 651 tion approach Zhang et al. (2024), the trajectories of all agents and the geometric representation of  
 652 lane segments are modeled as polylines composed of interconnected points. We employ an agent-  
 653 centric normalization strategy Cheng et al. (2023), which transforms all inputs into a coordinate  
 654 system centered on the target agent. The historical trajectories of  $N_a$  agents are represented as  $\mathbf{X}^n$ ,  
 655 which include coordinates and velocity changes over a length of  $T_{\max}$  timesteps. Furthermore, we  
 656 incorporate two time-related features, scaled timestamps to manage varying time ranges, and inter-  
 657 observation interval features to capture timing differences between observations, to construct the  
 658 time-aware input representation  $h^n$ .  
 659

660 Lane segments are encoded as  $\mathbf{h}^l$ , which includes the number of lane segments within a specified  
 661 radius around the target agent, the number of points in each polyline, and lane features such as  
 662 coordinates and availability. All coordinates within each lane segment are normalized relative to  
 663 their geometric centers, providing a standardized reference frame for subsequent processing and  
 664 analysis Bharilya et al. (2025).  
 665

666 **Evaluation Metric.** We assess our approach using several widely accepted metrics in trajectory  
 667 prediction research Wilson et al. (2021). The  $MinADE_k$  metric calculates the average Euclidean  
 668 distance between the predicted trajectories and the actual ground truth paths. The  $MinFDE_k$  met-  
 669 ric, on the other hand, measures the prediction error specifically at the endpoints of the trajectories.  
 670 To evaluate failure rates, the miss rate ( $MR_k$ ) counts instances where the endpoint error  $MinFDE_k$   
 671 exceeds a threshold of 2 meters. In these metrics,  $k$  denotes the number of trajectory modes pre-  
 672 dicted, with evaluations conducted for both single-mode predictions ( $k = 1$ ) and multi-modal pre-  
 673 dictions ( $k = 6$ ).  
 674

675 **Implementation Details.** The framework is implemented in PyTorch and trained on an NVIDIA  
 676 RTX A5000 GPU. Models are trained end-to-end for 60 epochs using the AdamW optimizer, with  
 677 a batch size of 128, a learning rate of 0.001, and a weight decay of 0.01. We use a cosine learning  
 678 rate schedule with a 10-epoch warm-up phase. An agent-centric coordinate system samples scene  
 679 elements within a 150-meter radius around the agents of interest. The embedding dimension is  
 680 set to 128. Each Mamba block contains 4 layers, the Transformer encoder has 5 layers, and the  
 681 Transformer decoder in the prediction module includes 2 layers.

682 **Formulation of Loss Functions.** To optimize the model, we use the Huber loss for trajectory  
 683 regression  $L_{reg}$  and cross-entropy loss for confidence classification  $L_{cls}$ . A winner-take-all strategy  
 684 is applied, optimizing only the best prediction while minimizing the average error relative to the  
 685 ground truth,  
 686

$$687 L_{reg} = \min_{k \in \{1, 2, \dots, 6\}} \left( \sum_{t=1}^{t_{pred}} \sum_{d=1}^2 L_{hl}(Y_{gt}^{(t,d)}, \hat{Y}_k^{(t,d)}) \right) \quad (24)$$

688 where  $\hat{Y}_k^{(t,d)}$  denotes the predicted future trajectory for mode  $k$  at timestamp  $t$  along coordinate  $d$ ,  
 689  $Y_{gt}^{(t,d)}$  represents the corresponding ground truth,  $t_{pred}$  is the total number of future time steps, and  $k$   
 690 indicates the number of predicted modes. For confidence classification, we apply the cross-entropy  
 691 loss,  
 692

$$693 L_{cls} = \sum_{k=1}^K \left( \mathbb{I}[Y_{gt}] \log(\pi_k) + (1 - \mathbb{I}[Y_{gt}]) \log(1 - \hat{\pi}_k) \right) \quad (25)$$

694 where  $\hat{\pi}_k$  is the predicted probability for the  $k$ -th trajectory, and  $\mathbb{I}[Y_{gt}]$  is an indicator function that  
 695 equals 1 if the  $k$ -th trajectory is closest to the ground truth, and 0 otherwise.  
 696

697 To supervise the missing-step temporal embedding, the loss  $L_{m-rg}$  is computed as,  
 698

$$700 L_{m-rg} = \sum_{m=1}^{t_m} \sum_{d=1}^2 L_{hl}(x_{gt}^{(m,d)}, \hat{x}_m^{(m,d)}) \quad (26)$$

701

where  $t_m$  is the number of missing data in the observation. Here,  $m$  indexes the missing points, and  $x_{\text{gt}}^{(m,d)}$  and  $\hat{x}_m^{(m,d)}$  denote the ground truth and predicted values of the missing data, respectively, along dimension  $d$ . The function  $L_{\text{hl}}$  refers to the Huber loss used for regression.

Moreover, to supervise the variable-step temporal embedding, the loss  $L_{v\text{-rg}}$  is defined as,

$$L_{v\text{-rg}} = \sum_{v=1}^{t_v} \sum_{d=1}^2 L_{\text{hl}} \left( x_{\text{gt}}^{(v,d)}, \hat{x}_v^{(v,d)} \right), \quad (27)$$

where  $t_v$  denotes the number of unobserved time steps with variable length. Here,  $v$  indexes the unobserved tokens, and  $x_{\text{gt}}^{(v,d)}$  and  $\hat{x}_v^{(v,d)}$  represent the ground truth and predicted values, respectively, along dimension  $d$ .

**Endpoint loss.** To predict endpoints, we utilize a dynamic multi-layer perceptron (MLP) with weights that are adaptively generated based on the input, referred to as the *adaptive MLP*. The adaptive MLP takes as input the agent features  $[E^t, F^n]$  and meta-information  $mi$  of all agents. The meta-information includes the agent’s position and normalized velocity at last observed timestamps. These inputs are concatenated and passed through an MLP with learnable parameters  $W_{\text{feat},1}$ ,  $W_{\text{feat},2}$  and biases  $b_{\text{feat},1}$ ,  $b_{\text{feat},2}$ , to obtain a latent representation,  $\tilde{f}$ ,

$$\tilde{f} = \varphi \left( W_{\text{feat},2} \varphi \left( W_{\text{feat},1} [E^t, F^n]; mi \right) + b_{\text{feat},1} \right) + b_{\text{feat},2} \quad (28)$$

with  $\varphi$  denoting the ReLU activation. Subsequently, two sets of dynamic weights  $W_1$  and  $W_2$  are generated by applying learnable linear transformations  $W_{d1}$  and  $W_{d2}$  to  $\tilde{f}$ , reshaped accordingly.

$$W_1 = \text{reshape}(W_{d1} \cdot \tilde{f} + b_{d1}) \quad (29)$$

$$W_2 = \text{reshape}(W_{d2} \cdot \tilde{f} + b_{d2}) \quad (30)$$

The first hidden layer activations  $F_{d1}$  are computed by applying a linear transformation  $W_1$  to the input feature  $f$ , followed by layer normalization and a non-linear activation function  $\varphi$ . The final prediction  $\hat{y}_{\text{ep}}$  is then obtained by applying a second linear transformation  $W_2$  to  $F_{d1}$ ,

$$F_{d1} = \varphi(\text{LayerNorm}(W_1 \cdot f)), \quad (31)$$

$$\hat{y}_{\text{ep}} = W_2 \cdot F_{d1} \quad (32)$$

This formulation enables dynamic adaptation of the prediction weights conditioned on the input features and meta information, allowing the model to flexibly predict agent endpoints. To improve the accuracy of endpoint predictions, we employ a dedicated loss defined as,

$$L_{\text{et}} = \sum_{n=1}^N \sum_{d=1}^2 L_{\text{hl}} \left( Y_{\text{gt},n}^{(t_{\text{pred}},d)}, \hat{e}p_n^{(t_{\text{pred}},d)} \right) \quad (33)$$

where  $L_{\text{et}}$  measures the discrepancy between the predicted endpoint  $\hat{e}p_n$  and the ground truth endpoint  $Y_{\text{gt},n}^{(t_{\text{pred}},d)}$  of the  $n^{\text{th}}$  agent, computed using the Huber loss function.

## A.2 MORE EXPERIMENTAL RESULTS

**Performance under Different Observation Settings.** The performance of the proposed model on the Argoverse 1 dataset under different input conditions is shown in Table 4.

**Ablation Study of BTD-Mamba Components.** Table 5 presents an ablation study analyzing different configurations of the BTD-Mamba module. Using only the forward Mamba (Fwd) or backward Mamba (Bwd) results in similar performance, with slightly better results for Fwd. Combining both directions (Fwd+Bwd) improves performance across all metrics, indicating that bidirectional context benefits trajectory modeling. The addition of Time Decay (TD) further enhances performance when combined with either Fwd or Bwd, showing that temporal relationships contribute useful dynamics. The best performance is achieved when all three components such as Fwd, Bwd, and TD are integrated, forming the complete BTD-Mamba module. This full configuration achieves the lowest minADE and minFDE, as well as the lowest miss rate, demonstrating the complementary nature of bidirectional processing and temporal differencing.

Table 4: Performance under different observation scenarios for Argoverse 1 dataset

Model	I/N Scenarios	minADE <sub>1</sub>	minFDE <sub>1</sub>	MR <sub>1</sub>	minADE <sub>6</sub>	minFDE <sub>6</sub>	MR <sub>6</sub>
DeMo-Orig	Missing + Var.	2.3924	4.6935	0.6172	0.9052	1.4566	0.1751
	Var. Obs	2.1506	4.2475	0.5993	0.8311	1.3058	0.5993
	Missing Only	3.4942	6.6470	0.7808	1.4078	2.2559	0.3238
	Full Obs	1.2903	2.7863	0.46495	0.5926	0.9534	0.0830
DeMo-RSD	Missing + Var.	1.5100	3.1800	0.5174	0.6504	1.0569	0.1010
	Var. Obs	1.5798	3.3039	0.5380	0.6538	1.0576	0.1022
	Missing Only	2.7708	1.8213	0.7421	1.0927	1.8213	0.2572
	Full Obs	1.3999	2.9870	0.5030	0.6145	0.9946	0.0918
Forecast-mae-Orig	Missing + Var.	1.7647	3.6791	0.6153	0.7315	1.2033	0.1175
	Var. Obs	1.4734	3.1700	0.5340	0.6669	1.0936	0.0964
	Missing Only	2.2504	4.4974	0.6836	0.8654	1.4303	0.1757
	Full Obs	1.3470	2.9207	0.0901	0.6223	1.0222	0.0901
Forecast-mae-RSD	Missing + Var.	1.4679	3.1545	0.5226	0.6643	1.0849	0.0952
	Var. Obs	1.4576	3.1298	0.5243	0.6599	1.0759	0.0945
	Missing Only	1.7459	3.6371	0.6066	0.7206	1.1881	0.1172
	Full Obs	1.4338	3.0894	0.5120	0.6554	1.0707	0.0949
Our	Missing + Var.	0.8212	1.7826	0.3028	0.3956	0.6018	0.0698
	Var. Obs	0.8086	1.7343	0.2871	0.3693	0.5901	0.0510
	Missing Only	0.8123	1.7365	0.2816	0.3756	0.5902	0.0511
	Full Obs	0.7943	1.7029	0.2718	0.3346	0.5726	0.0467

**Ablation Study on Predictor Module Components.** Table 6 evaluates the impact of three components in the predictor module: State Encoding for Intermittent Trajectories (SEIT), Multi-Modal Trajectory Distributions (MMTD), and Unified Mode Fusion (UMF). When MMTD is not used, separate MLPs replace it for mode and state prediction. The baseline model without these components (ID-1) performs the worst. Incorporating each component individually (ID-2 to ID-4) yields moderate improvements, indicating their standalone effectiveness. Combinations of two components (ID-5 to ID-7) further enhance performance, demonstrating complementary strengths. The full model with all three components enabled (ID-8) achieves the best results, confirming that SEIT, MMTD, and UMF together significantly improve trajectory prediction accuracy across all metrics.

**Ablation Study on the Depth of BTM-Mamba and Transformer Encoder.** Table 7 investigates the impact of varying the depth i.e., the number of stacked layers, of the BTM-Mamba module and the Transformer Encoder on model performance. Increasing the depth of both modules generally improves results, as seen when moving from 3 to 4 layers, leading to reduced minADE, minFDE, and Miss Rate (MR). The best performance is observed with 4 layers of BTM-Mamba and 5 layers of the Transformer Encoder, achieving the lowest across all metrics. Moreover, performance slightly declines when the BTM-Mamba depth is increased to 5 layers alongside 5 Transformer layers, suggesting a trade-off where excessive depth in BTM-Mamba fails to yield further benefits. Overall, a moderate depth configuration balances model complexity and predictive accuracy effectively.

**Ablation Study on the Impact of Auxiliary Losses.** Table 8 presents an ablation study that evaluate contribution of three auxiliary losses: the endpoint loss ( $L_{et}$ ), regression loss over variable-step temporal embedding ( $L_{u-r_g}$ ), and regression loss for missing-step temporal embedding ( $L_{m-r_g}$ ). The baseline model without any auxiliary loss (ID-1) shows the weakest performance across all metrics. Introducing each loss individually (ID-2 to ID-4) yields consistent improvements, demonstrating their individual effectiveness. Both  $L_{u-r_g}$  and  $L_{m-r_g}$  lead to greater gains than  $L_{et}$ . Combining two of the losses (ID-5 to ID-7) further improves performance, showing their complementary effects. The best results are obtained when all three auxiliary losses are applied simultaneously (ID-8), achieving the lowest minADE, minFDE, and MR. These findings confirm that auxiliary supervision strengthens the model’s ability to learn more accurate trajectory representations.

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

Table 5: Component Study of BTM-Mamba for Argoverse 2 dataset

Method	minADE <sub>1</sub>	minFDE <sub>1</sub>	MR <sub>1</sub>	minADE <sub>6</sub>	minFDE <sub>6</sub>	MR <sub>6</sub>
Fwd	1.6816	4.1148	0.6027	0.7011	1.3463	0.1804
Bwd	1.6901	4.1757	0.6134	0.7128	1.3524	0.1749
Fwd+Bwd	1.6529	4.0284	0.5903	0.6925	1.3163	0.1628
Fwd+TD	1.6425	4.0143	0.5901	0.6911	1.3047	0.1609
Bwd+TD	1.6546	4.0112	0.5928	0.6836	1.3142	0.1628
BTM-Mamba	1.5882	3.9402	0.5717	0.6562	1.2326	0.1551

Table 6: Component Study of Predictor Module for Argoverse 2 dataset

ID	SEIT	MMTD	UMF	minADE <sub>1</sub>	minFDE <sub>1</sub>	MR <sub>1</sub>	minADE <sub>6</sub>	minFDE <sub>6</sub>	MR <sub>6</sub>
1	-	-	-	1.6474	4.2143	0.6245	0.6957	1.2764	0.1802
2	✓	-	-	1.6349	4.2094	0.6137	0.6835	1.2638	0.1782
3	-	✓	-	1.6298	4.2072	0.6048	0.6804	1.2477	0.1756
4	-	-	✓	1.6164	4.1537	0.6013	0.6787	1.2416	0.1726
5	✓	✓	-	1.6064	4.1121	0.5912	0.6765	1.2569	0.1683
6	✓	-	✓	1.6092	4.0154	0.5936	0.6627	1.2535	0.1647
7	-	✓	✓	1.5914	3.9668	0.5805	0.6613	1.2476	0.1589
8	✓	✓	✓	1.5882	3.9402	0.5717	0.6562	1.2326	0.1551

Table 7: Depth Study of BTM-Mamba and Transformer Encoder for Argoverse 2 dataset

	BTM-M	T-Enc	minADE <sub>1</sub>	minFDE <sub>1</sub>	MR <sub>1</sub>	minADE <sub>6</sub>	minFDE <sub>6</sub>	MR <sub>6</sub>
	3	3	1.6226	4.0151	0.5926	0.6844	1.2531	0.1629
	4	4	1.6048	3.9197	0.5802	0.6614	1.2494	0.1600
	4	5	1.5882	3.9402	0.5717	0.6562	1.2326	0.1551
	5	5	1.5901	3.9937	0.5826	0.6632	1.2471	0.1622

Table 8: Impact of Auxiliary Losses for Argoverse 2 dataset

ID	$L_{et}$	$L_{u-rg}$	$L_{m-rg}$	minADE <sub>1</sub>	minFDE <sub>1</sub>	MR <sub>1</sub>	minADE <sub>6</sub>	minFDE <sub>6</sub>	MR <sub>6</sub>
1	-	-	-	1.6284	4.2591	0.6137	0.6814	1.2641	0.1795
2	✓	-	-	1.6211	4.1918	0.6093	0.6787	1.2601	0.1725
3	-	✓	-	1.6159	4.0166	0.5935	0.6749	1.2538	0.1617
4	-	-	✓	1.6117	4.0137	0.5931	0.6732	1.2546	0.1629
5	✓	✓	-	1.6026	4.0118	0.5874	0.6726	1.2525	0.1658
6	✓	-	✓	1.5984	3.9971	0.5846	0.6815	1.2437	0.1604
7	-	✓	✓	1.5907	3.9473	0.5824	0.6810	1.2429	0.1598
8	✓	✓	✓	1.5882	3.9402	0.5717	0.6562	1.2326	0.1551



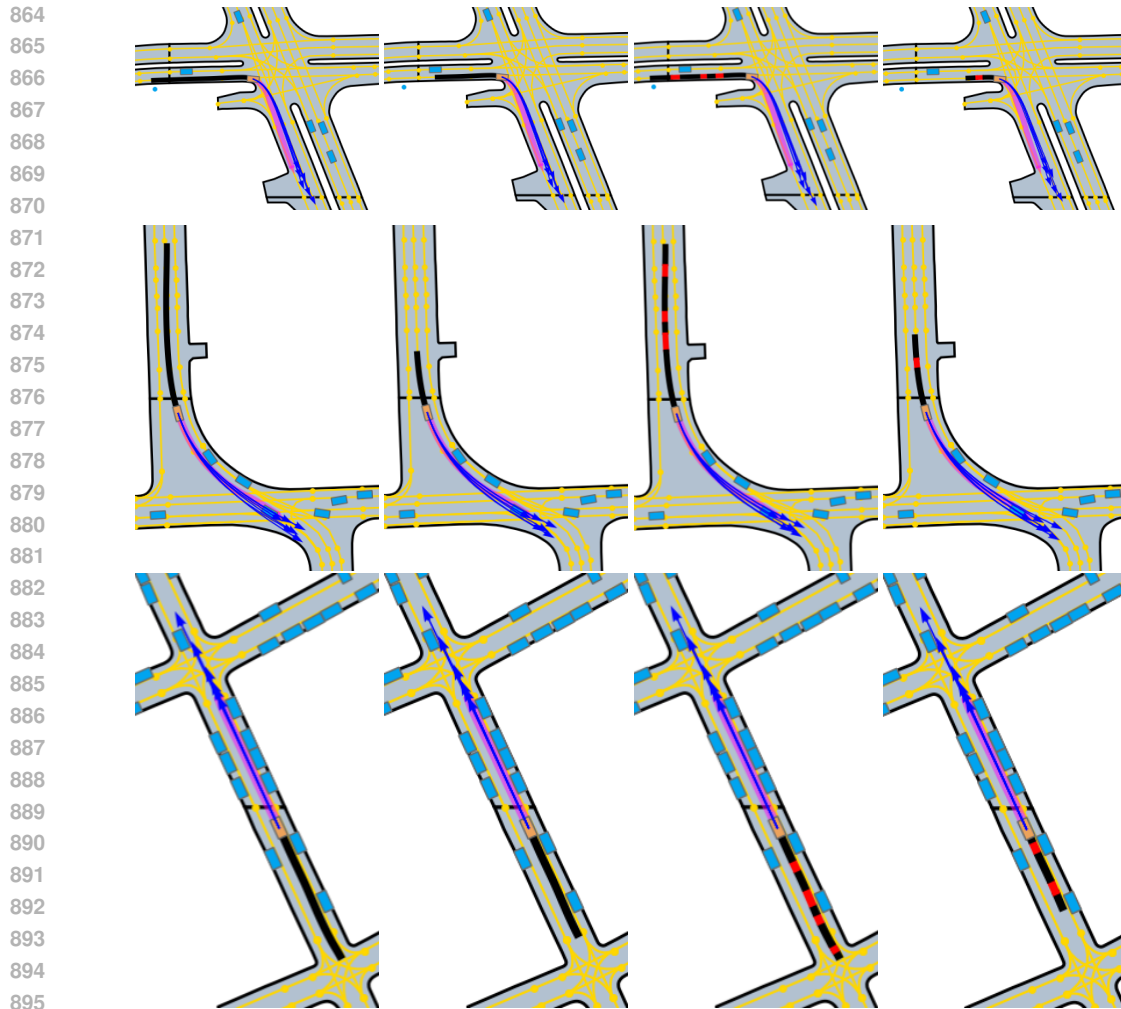


Figure 4: Qualitative results of the proposed model with varying input observations on Argoverse 2 dataset. Black: observed history; red: missing points; magenta: ground truth future; blue: predicted future trajectory.

### A.3 MORE QUALITATIVE RESULTS

The qualitative results are presented in Figure 4, showcasing the performance of the proposed model across diverse driving scenarios. The first row illustrates an intersection scenario where the agent executes a right turn, requiring awareness of both lane geometry and surrounding context. The second row demonstrates the agent’s behavior on a curved path, highlighting the model’s ability to capture smooth trajectory changes over time. The third row presents a straight-driving scenario in dense traffic, where the model must accurately predict future motion despite limited maneuvering space and potential occlusions. Across all scenarios, the columns depict different input conditions, including variable-length observations and missing data. The proposed model consistently produces coherent and accurate trajectory predictions, demonstrating its adaptability to a wide range of real-world input conditions.