



# Deep Learning-based Trajectory Prediction using Improved HiVT

System Control Lab.

44211032-5 SUN Honglin

2022/2/6

# Contents

---

## 1. Introduction

~~1.1 Background & Research Target~~

1.2 Issue & Idea

---

## 2. Deep Learning-based Trajectory Prediction

2.0 Intro to HiVT (Hierarchical Vector Transformers)

2.1 Computational Efficiency Improvement

2.2 Prediction Accuracy Improvement

---

## 3. Simulations

3.1 Simulation Setup

3.3 Quantitative Comparisons

3.2 Visualization Results

3.4 Effectiveness Validation

---

## 4. Conclusions



## 1.2 Issue & Idea

### Issue

- ❑ Most learning-based methods require **heavy computation**  
→ Hard to apply to vehicles
- ❑ Long **latency** will make the prediction meaningless

### Idea

Perform a two-stage modification on **HiVT-64**:

- ❑ Improve the **computational efficiency**
  - ~~Cut down parameters and layers~~
  - ~~Modify or remove some operations~~
  - Implement a new local encoder
- ❑ Improve the prediction **accuracy**
  - Refine the network structure
  - Introduce more useful features
  - ~~Improve the training strategy~~

TABLE I: Comparison of some deep learning-based methods (LaneGCN as the baseline of accuracy)

Method	Accuracy	#Param	Device (GPU)
LaneGCN	100.00%	3,701K	GTX TITAN X
Scene Transformer	110.71%	15,296K	Tesla V100
GOHOME	94.05%	400K	RTX 2080 Ti
<b>HiVT-64</b> [4]	99.77%	653K	RTX 2080 Ti

TABLE II: Comparison of some devices (RTX 2080 Ti as the baseline of performance)

Device	Performance	Power
RTX 2080 Ti	100.0%	250 W
<b>GTX 1070 Max-Q</b>	42.0%	115 W
Jetson AGX Xavier	10.5%	30 W
Jetson AGX Orin	39.6%	60 W

Onboard  
AI chips



[4] Zhou, Zikang, et al. "HiVT: Hierarchical Vector Transformer for Multi-Agent Motion Prediction." *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022.

# Contents

---

## 1. Introduction

~~1.1 Background & Research Target~~

1.2 Issue & Idea

---

## 2. Deep Learning-based Trajectory Prediction

2.0 Intro to HiVT (Hierarchical Vector Transformers)

2.1 Computational Efficiency Improvement

2.2 Prediction Accuracy Improvement

---

## 3. Simulations

3.1 Simulation Setup

3.3 Quantitative Comparisons

3.2 Visualization Results

3.4 Effectiveness Validation

---

## 4. Conclusions



# 2.0 Intro to HiVT (Hierarchical Vector Transformers)

## Local regions

- **Small circular regions** centered at each agent

## Modules inside the model (4 stacked Transformers)

- Agent-Agent interaction  $\times T$   
Encode **agents' interaction** in **local region** at each time step
- Temporal transformer  
Capture **temporal information**
- Agent-Lane interaction  
Encode the **map information**
- Global interaction  $\times 3$   
Encode the **relationship between local regions**  
(Perform 3 times)

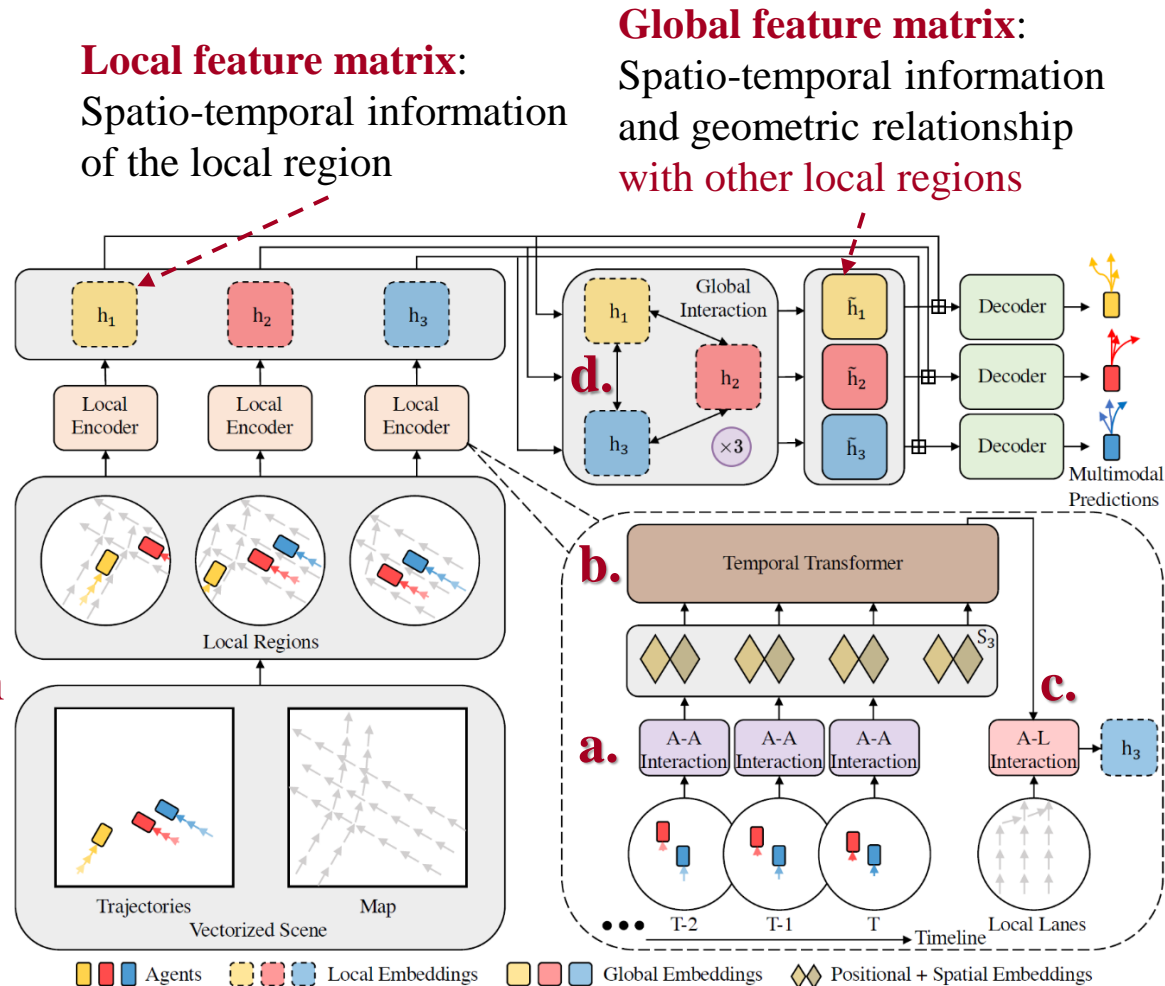


Fig. 4: Overview of HiVT



## 2.1 Computational Efficiency Improvement

### Implement a new local encoder

#### □ Original local encoder

1. Agent-Agent Interaction  $\times T$

2. Temporal Transformer

3. Agent-Lane Interaction

□ A-A interaction module:  
time-consuming

■ Not all agents have interaction  
with each other

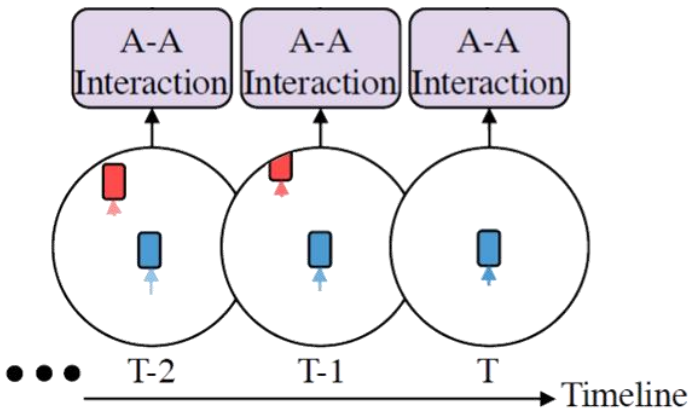


Fig. 13: Issue of A-A interaction

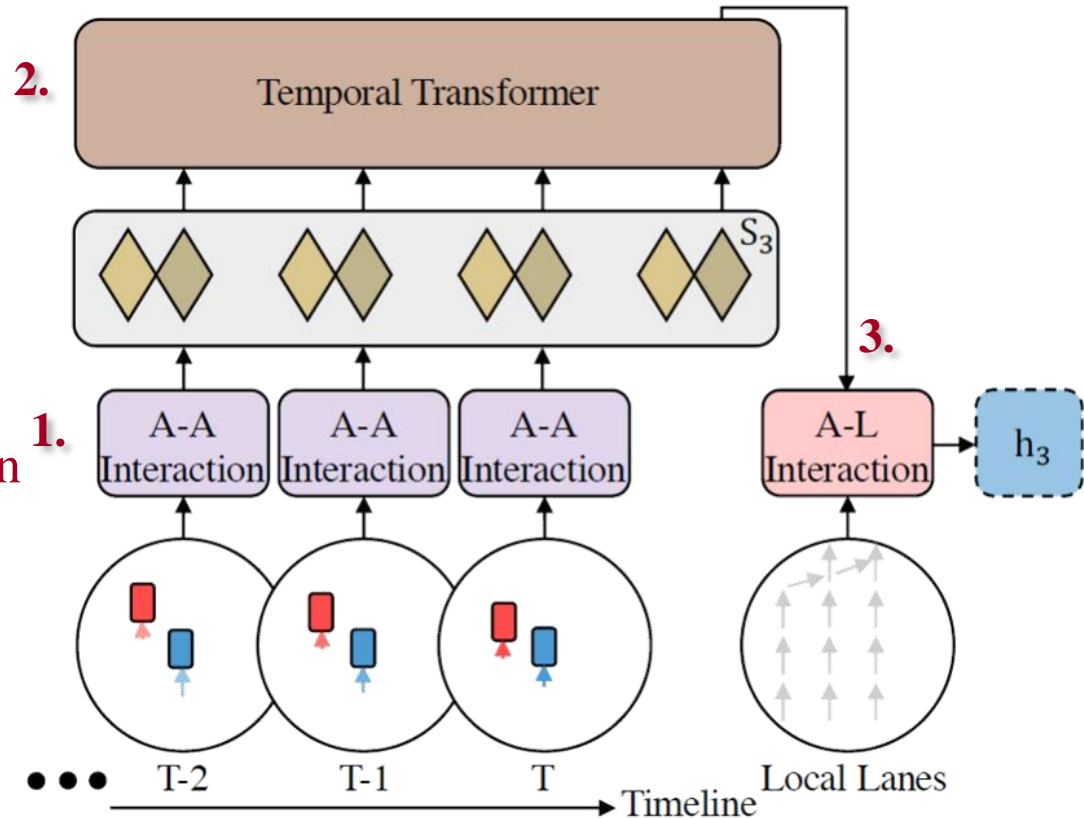


Fig. 12: Original local encoder



# 2.1 Computational Efficiency Improvement

## Implement a new local encoder

### □ Original local encoder

1. Agent-Agent Interaction  $\times T$
2. Temporal Transformer
3. Agent-Lane Interaction

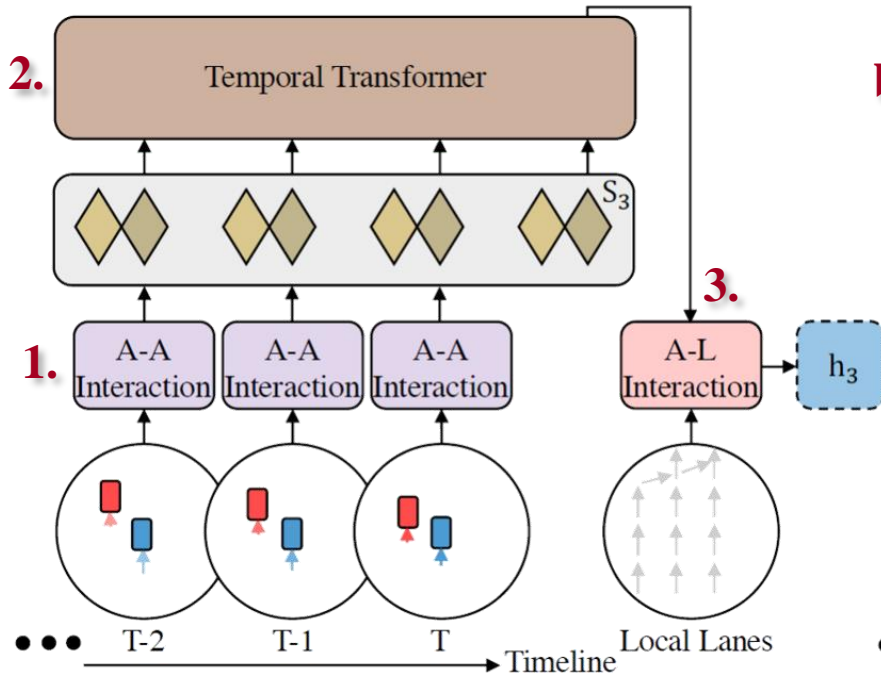


Fig. 12: Original local encoder

### □ Proposed local encoder

- a. Agent Encoder  $\times T$
  - b. Temporal Transformer
  - c. Agent-Lane Interaction
  - d. Agent-Agent Interaction  $\times 1$
- A-L before A-A:**  
Enrich the features of each agent

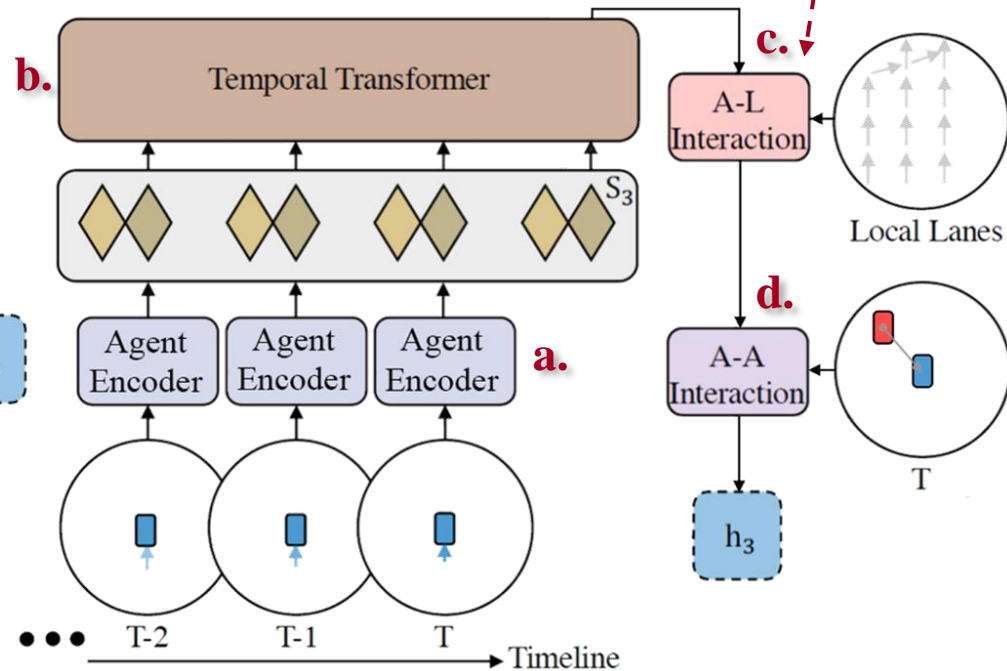


Fig. 14: Proposed local encoder



## 2.2 Prediction Accuracy Improvement

### Refine the network structure

#### □ Modify the **trajectory generation**

##### ■ Change the **decoder** output size:

30-position **dense** output ( $0.1\text{s} \times 30$ )

→ 3-position **sparse** output ( $1\text{s} \times 3$ )

□ Fewer data need to learn:  $30 \times 2 \rightarrow 3 \times 2$

□ Makes the model **easier to converge**

##### ■ Use cubic **B-spline interpolation** to complete the full trajectory

□ **Driving purpose** with 1-second interval is almost **deterministic**

→ Enough for the mathematical method to interpolate an accurate trajectory

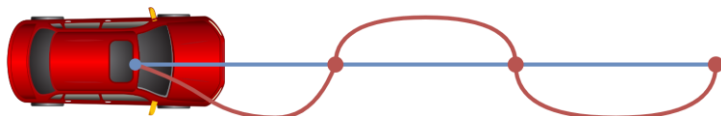
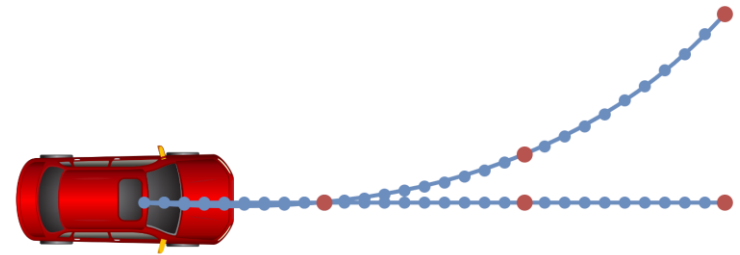
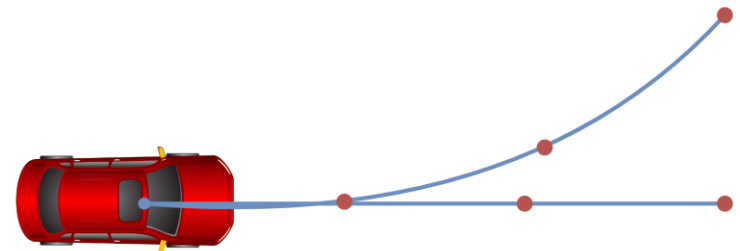


Fig. 16: Example of a driving purpose (go straight)

●: Every 1 second  
●: Every 0.1 second



(a) Directly output all positions (original)



(b) Sparse output with interpolation (proposed)

Fig. 15: Comparison of different output





## 2.2 Prediction Accuracy Improvement

### Refine the network structure

□ Modify the **fusion method of edge attributes** ( $e$ ) inside global interaction module

#### ■ Original

$$\mathbf{e}_{ij} = \phi_{\text{MLP}}([R_{ij}\mathbf{p}_{ij}, \sin(\theta_{ij}), \cos(\theta_{ij})])$$

$$\mathbf{h}_j = \mathbf{h}_j$$

$$\rightarrow Q = W_{Q_{\text{node}}} \mathbf{h}_i$$

$$K = W_{K_{\text{node}}} \mathbf{h}_j + W_{K_{\text{edge}}} \mathbf{e}_{ij}$$

$$V = W_{V_{\text{node}}} \mathbf{h}_j + W_{V_{\text{edge}}} \mathbf{e}_{ij}$$

#### ■ Proposed

$$\mathbf{e}_{ij} = [R_{ij}\mathbf{p}_{ij}, \sin(\theta_{ij}), \cos(\theta_{ij})]$$

$$\mathbf{h}_j = \phi_{\text{MLP}}([\mathbf{h}_j, \mathbf{e}_{ij}])$$

$$\rightarrow Q = W_{Q_{\text{node}}} \mathbf{h}_i$$

$$K = W_{K_{\text{node}}} \mathbf{h}_j \leftarrow \text{Retains more node information}$$

$$V = W_{V_{\text{node}}} \mathbf{h}_j$$

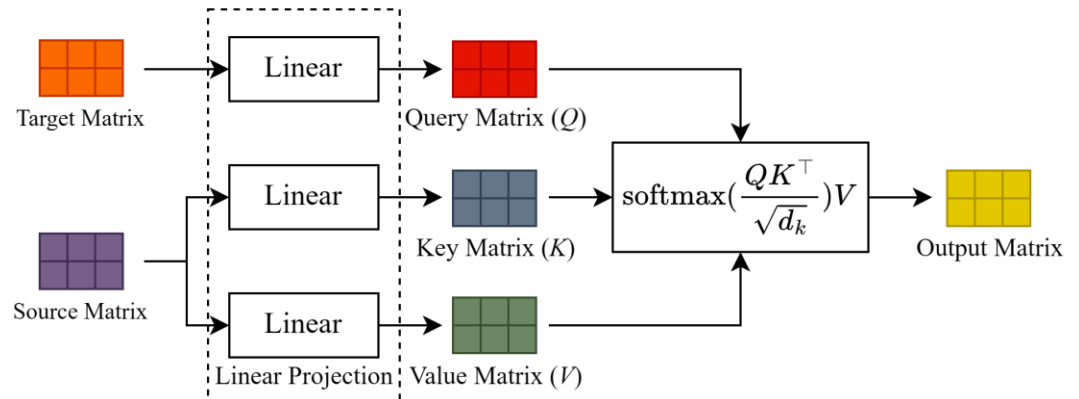
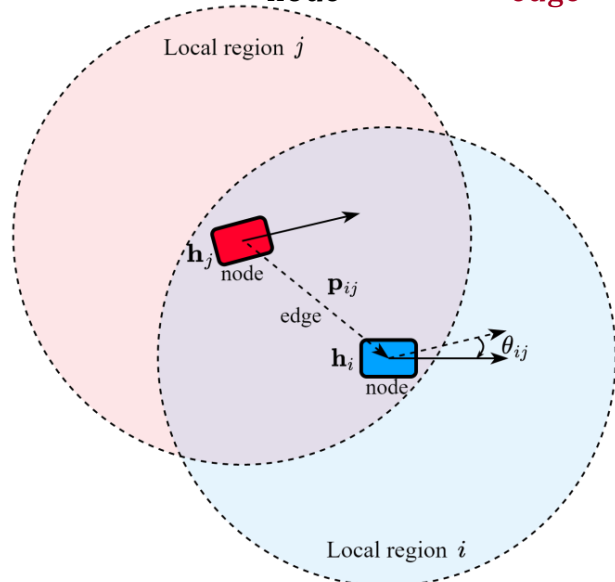


Fig. 7: Scaled-dot product attention inside the transformers

Fig. 17: Illustration of edge attributes



## 2.2 Prediction Accuracy Improvement

### Introduce more useful features

- Add learnable **time weights** to the temporal transformer

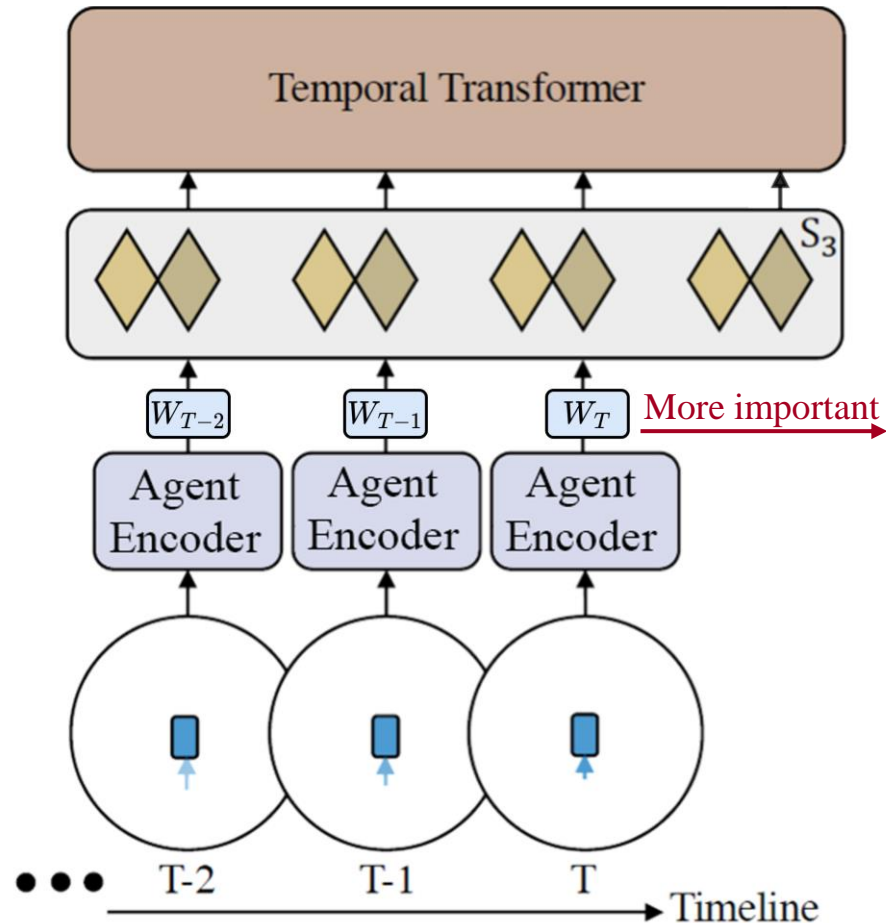


Fig. 21: Time weights

- Fusing **lane adjacency relationship** as semantic information

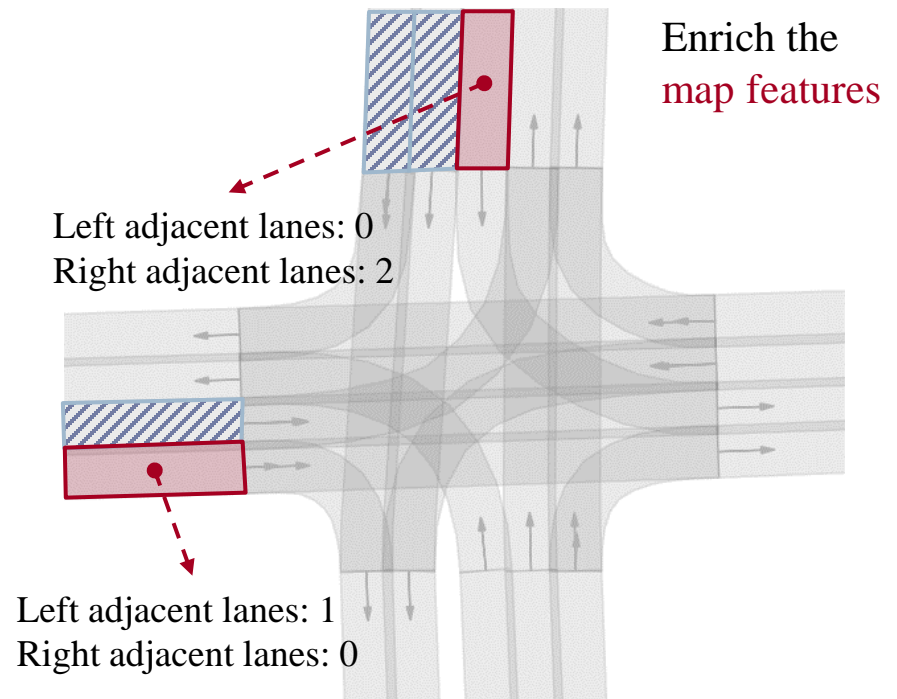


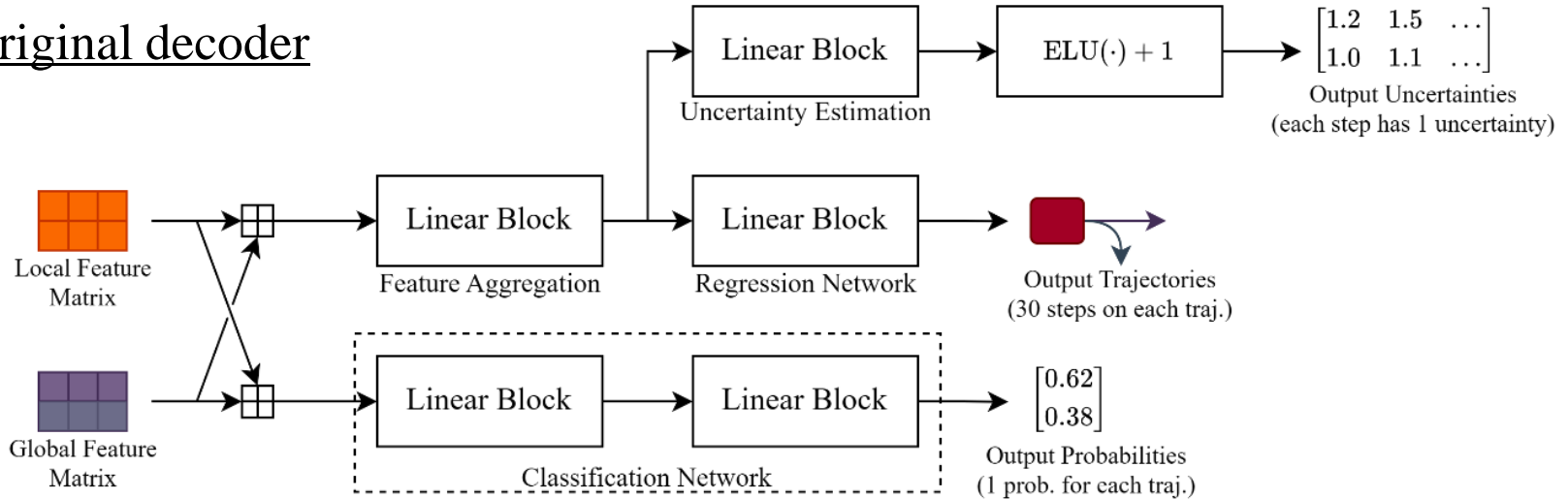
Fig. 22: Example usage of lane adjacency relationship



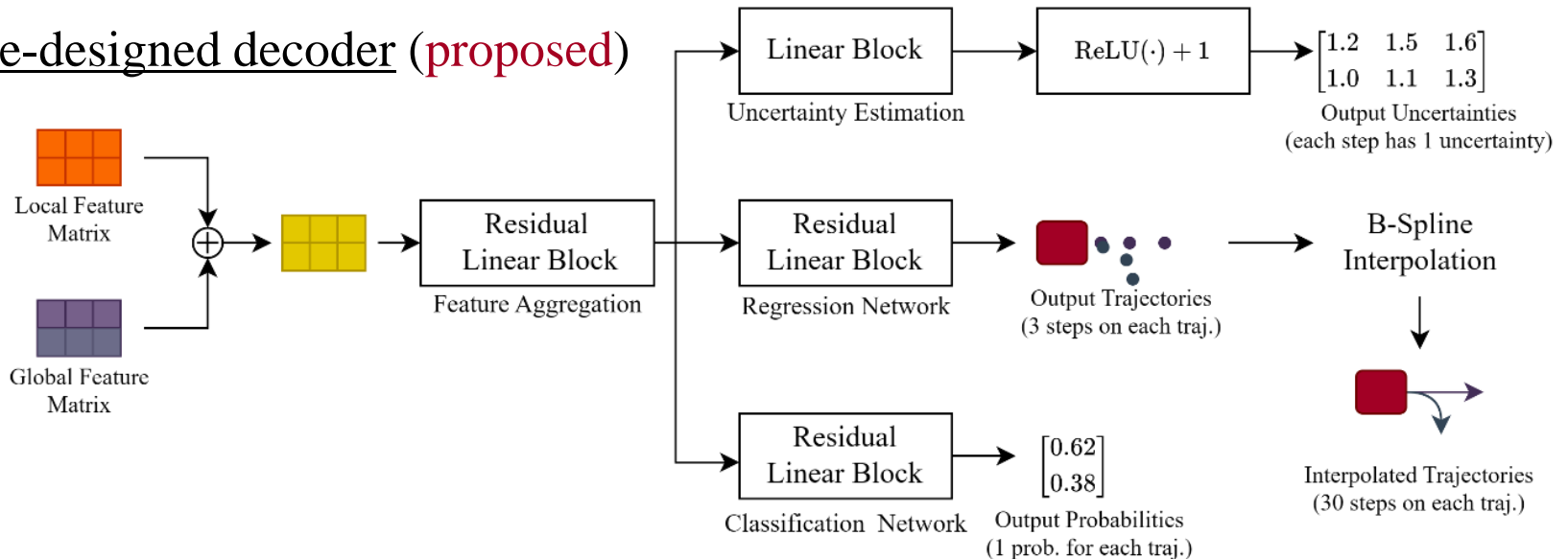
## 2.2 Prediction Accuracy Improvement

### Overall comparison of the decoder structure

#### Original decoder

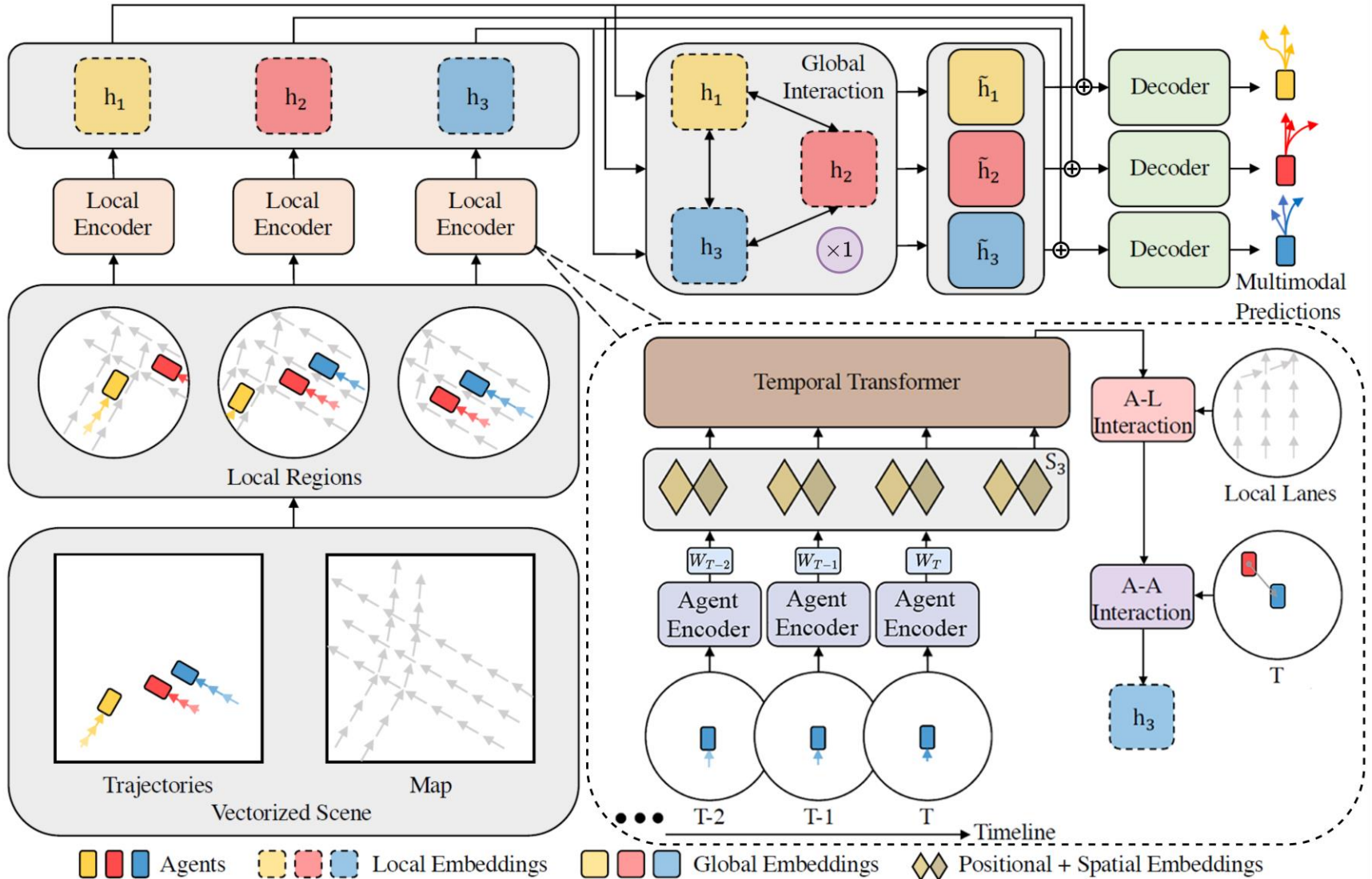


#### Re-designed decoder (proposed)



## 2.2 Prediction Accuracy Improvement

### Overall structure of the modified HiVT



# Contents

---

## 1. Introduction

~~1.1 Background & Research Target~~

1.2 Issue & Idea

---

## 2. Deep Learning-based Trajectory Prediction

2.0 Intro to HiVT (Hierarchical Vector Transformers)

2.1 Computational Efficiency Improvement

2.2 Prediction Accuracy Improvement

---

## 3. Simulations

3.1 Simulation Setup

3.3 Quantitative Comparisons

3.2 Visualization Results

3.4 Effectiveness Validation

---

## 4. Conclusions



## 3.1 Simulation Setup

### Argoverse 1 Motion Forecasting Dataset [6]

- ❑ The same as used in HiVT and other papers
- ❑ Data was collected in Miami and Pittsburgh
- ❑ 324,557 scenarios in the dataset, split into:
  - Training set: 205,942 scenarios
  - Validation set: 39,472 scenarios
  - Test set: 78,143 scenarios
- ❑ 5 seconds of data for each scenario
  - 10 Hz sampling rate
- ❑ Task
  - 2-second historical trajectory + Map
  - ↓ Predict
  - 3-second possible future trajectories  
(30 positions  $\times$  number of trajs.)

```
lanes[lane_id] = {  
    centerline = [X,Y]  
    is_intersection = True  
    turn_direction = 'LEFT'  
    traffic_controls = True  
    left_neighbor = None  
    right_neighbor = lane_id  
    predecessors = [lane_ids]  
    successors = [lane_ids]  
}
```

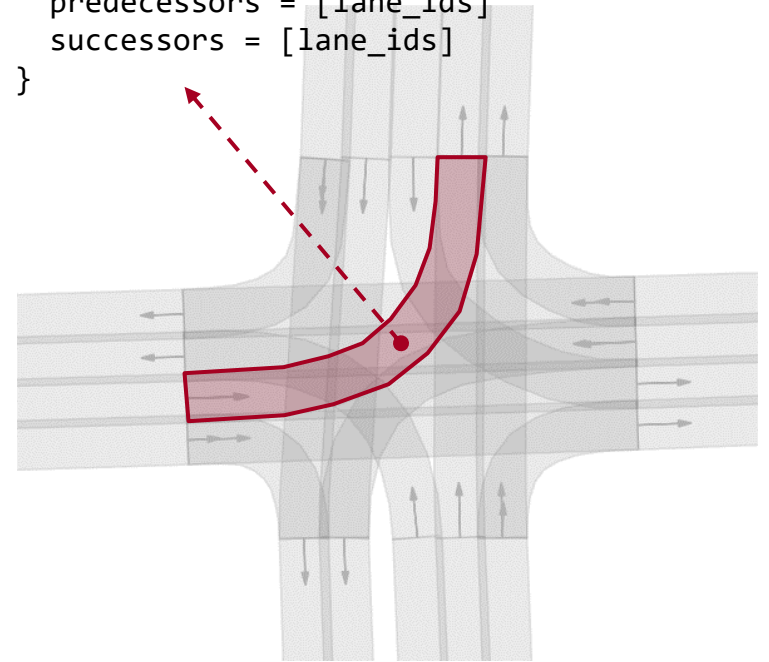


Fig. 24: Example of map data structure





# 3.1 Simulation Setup

## Visualization

- Predict 6 possible future trajectories
  - Yellow line: Historical trajectory (2s)
  - Red line: True future trajectory (3s)
  - Green lines: Predicted trajectories, probability declines from 1 to 6
  - Cyan dots: Final positions of other agents
  - Gray polygons: Lane segments
  - Gray arrows: Lane directions
- Good prediction: Close to the red line
- Plot one vehicle's prediction result
- Darker green indicates a higher probability

K=1: Consider the highest prob. trajectory

K=6: Consider all the 6 trajectories

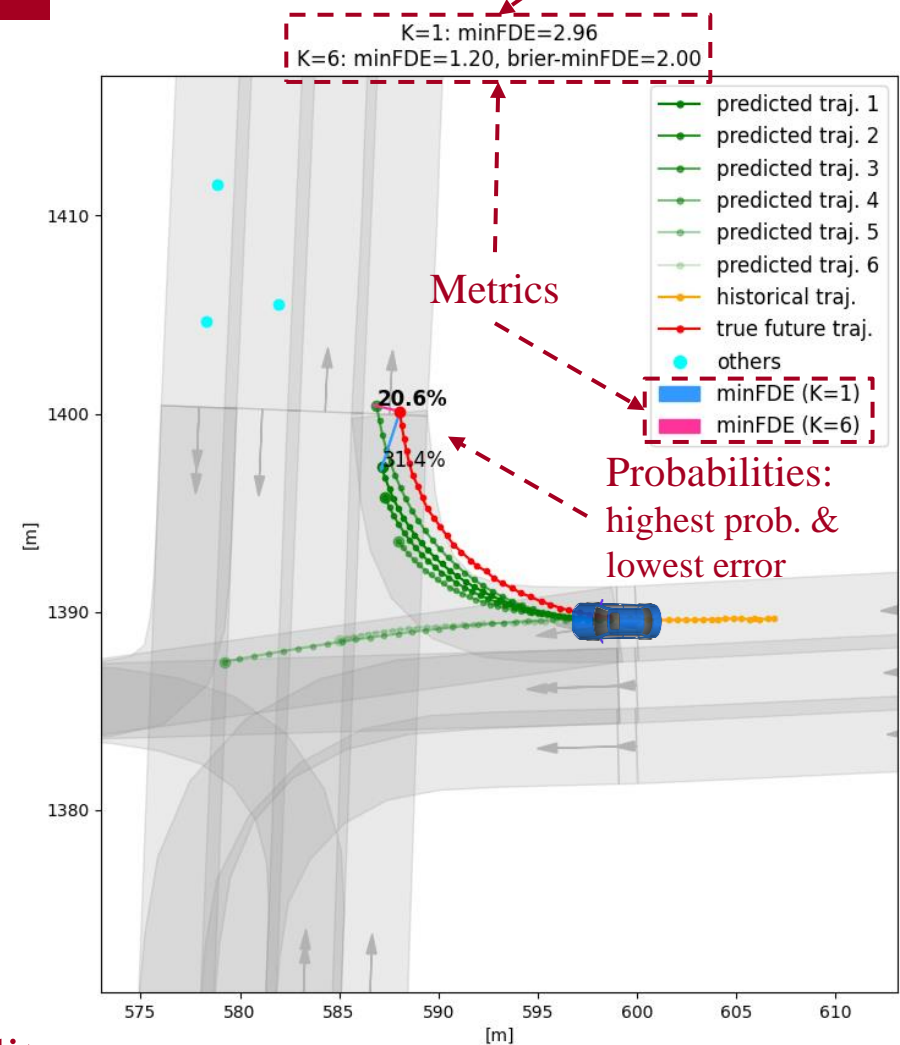


Fig. 25: Example of prediction result



WASEDA UNIVERSITY



# 3.1 Simulation Setup

**Metrics** (Lower is better)

□ minADE: Minimum **Average** Displacement Error

$$\text{minADE} = \min \left( \frac{1}{T} \sum_{t=1}^T \|\hat{\mathbf{y}}_t - \mathbf{y}_t\|_2 \right)$$

□ minFDE: Minimum **Final** Displacement Error

$$\text{minFDE} = \min(\|\hat{\mathbf{y}}_T - \mathbf{y}_T\|_2)$$

□ brier-minFDE [7]: minFDE + **probability term**  
→ Evaluate both **regression** and **classification** performance

$$\text{brier-minFDE} = \min(\|\hat{\mathbf{y}}_T - \mathbf{y}_T\|_2 + (1 - p)^2)$$

□ MR: Miss Rate

$$\text{MR} = \frac{\#\text{scenarios\_with\_minFDE} > 2[\text{m}]}{\#\text{scenarios}}$$

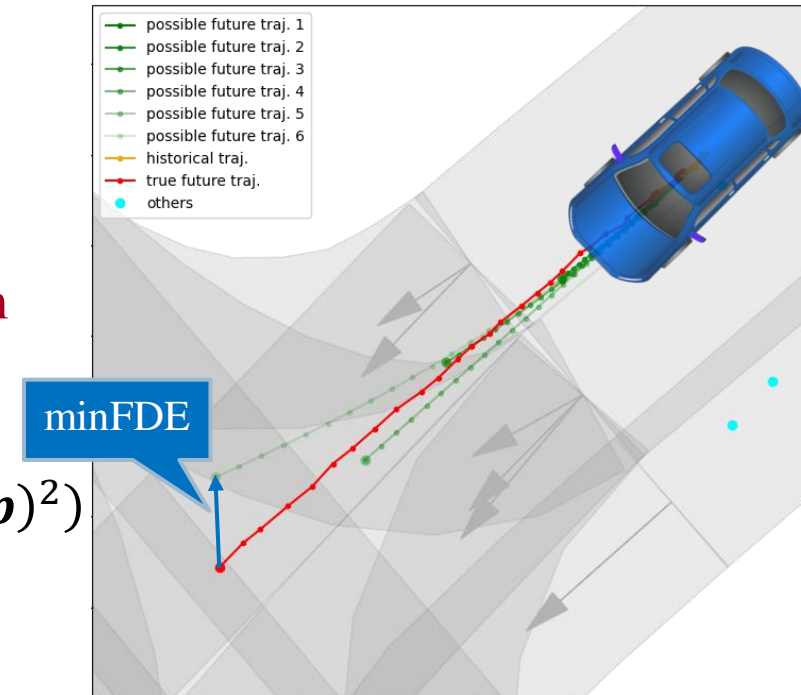


Fig. 29: Calculation of minFDE

$\hat{\mathbf{y}}_t$ : Predicted future **positions** at time step  $t$

$\mathbf{y}_t$ : True position at time step  $t$

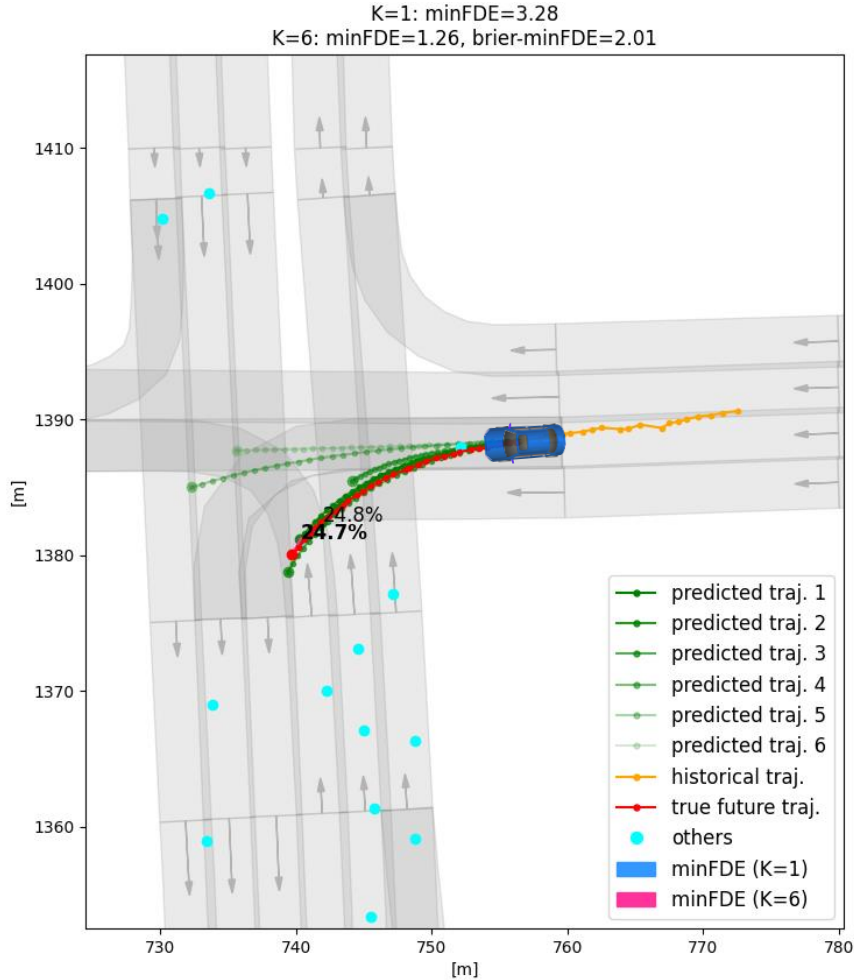
$p$ : Predicted **probabilities** of future positions

[7] Brier, Glenn W. "Verification of Forecasts Expressed in Terms of Probability." *Monthly Weather Review* 78.1 (1950): 1-3.

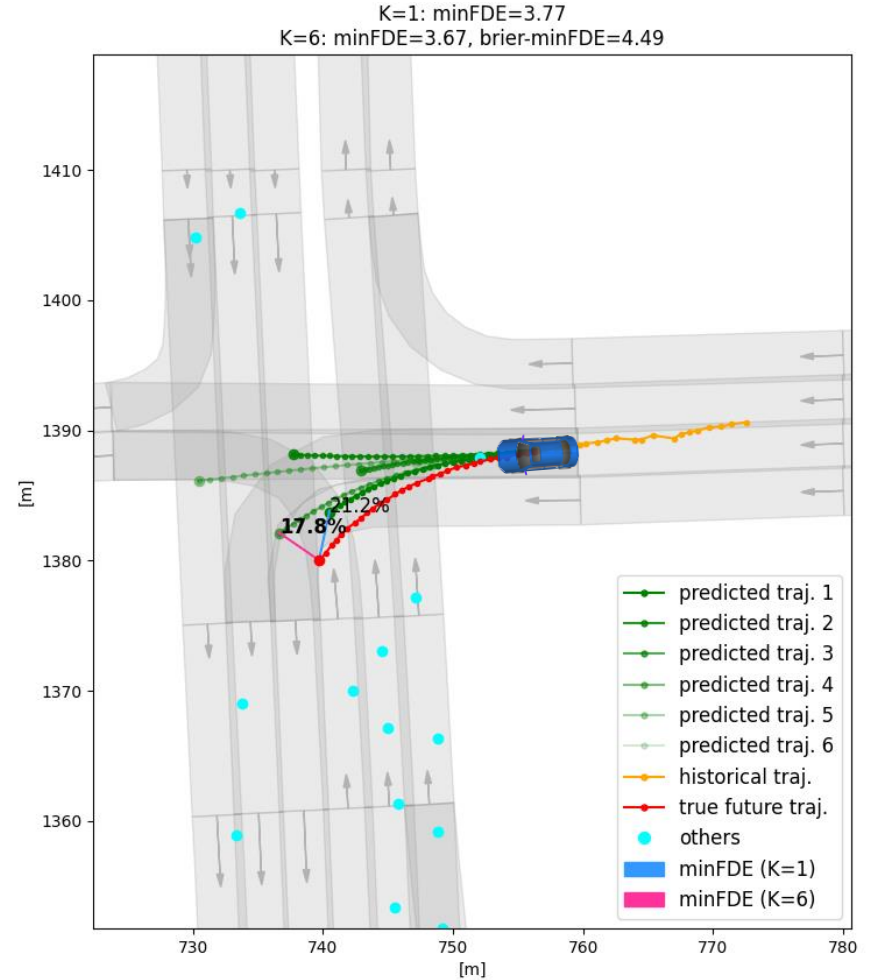


## 3.2 Visualization Results

### Comparison with the original HiVT-64 model



(a) Proposed

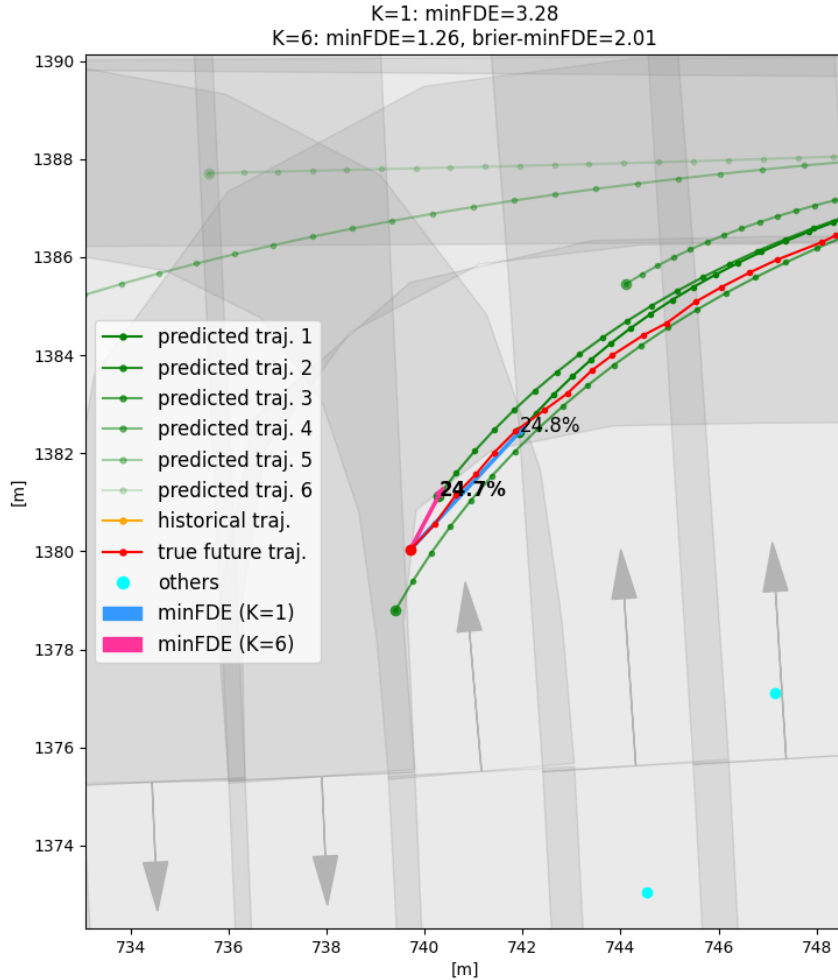


(b) Original

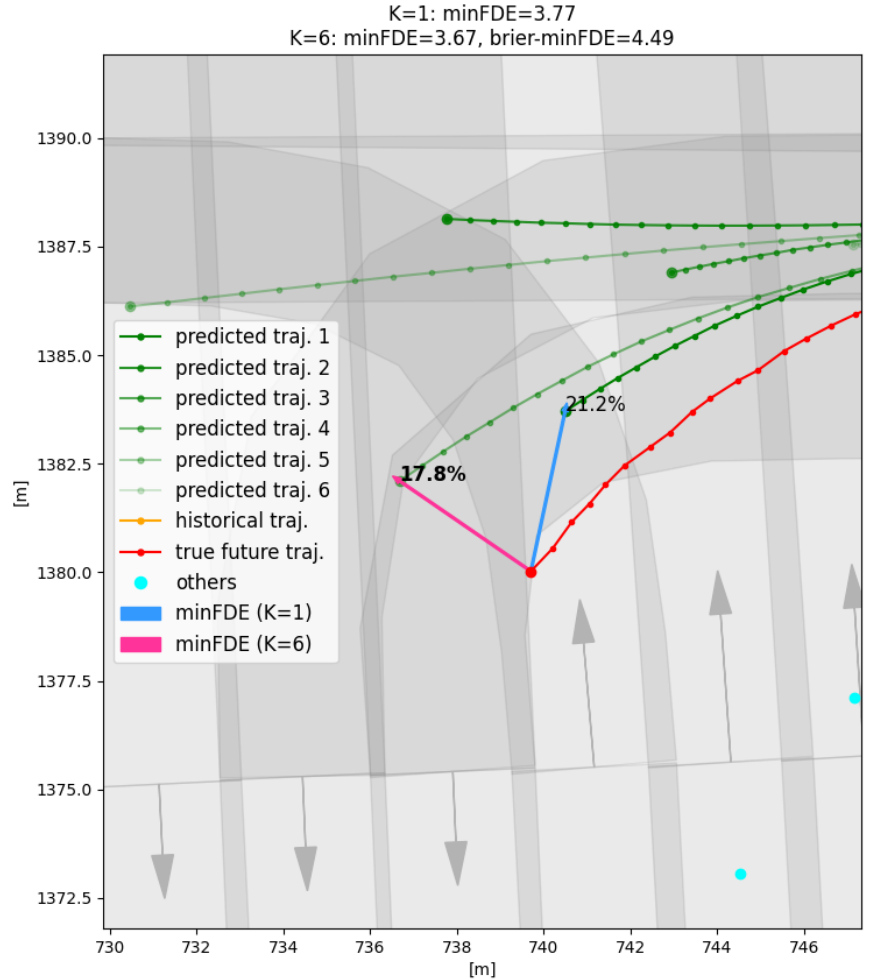
Fig. 26: Turning left

## 3.2 Visualization Results

### Comparison with the original HiVT-64 model (zoom in)



(a) Proposed

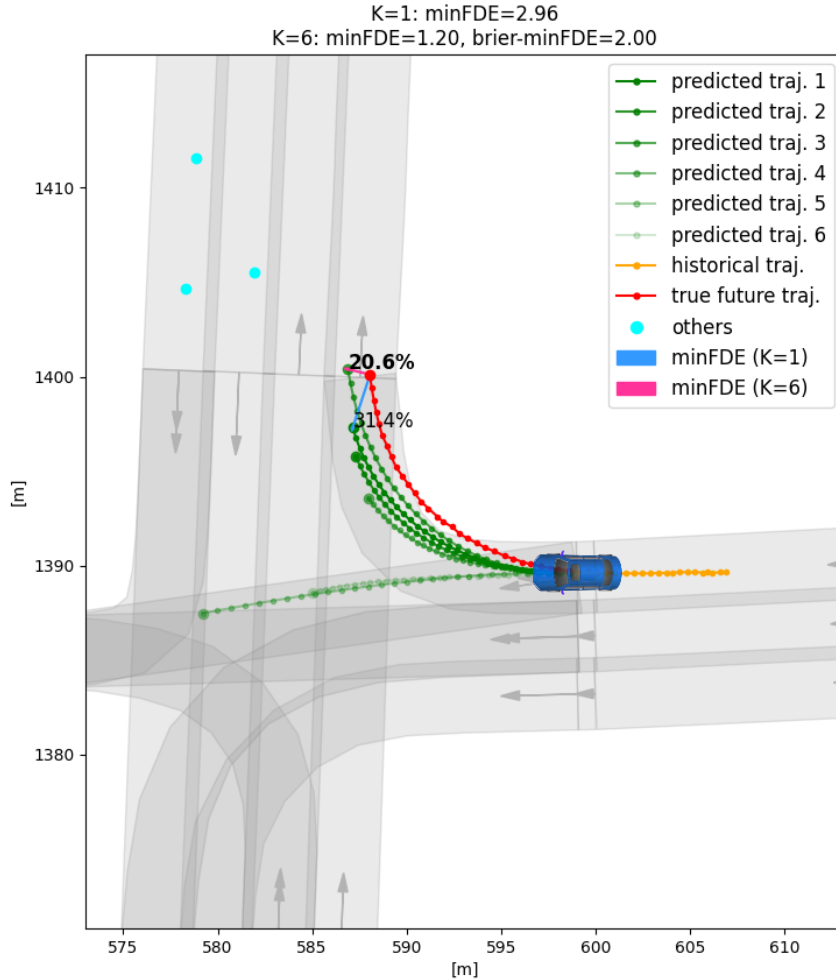


(b) Original

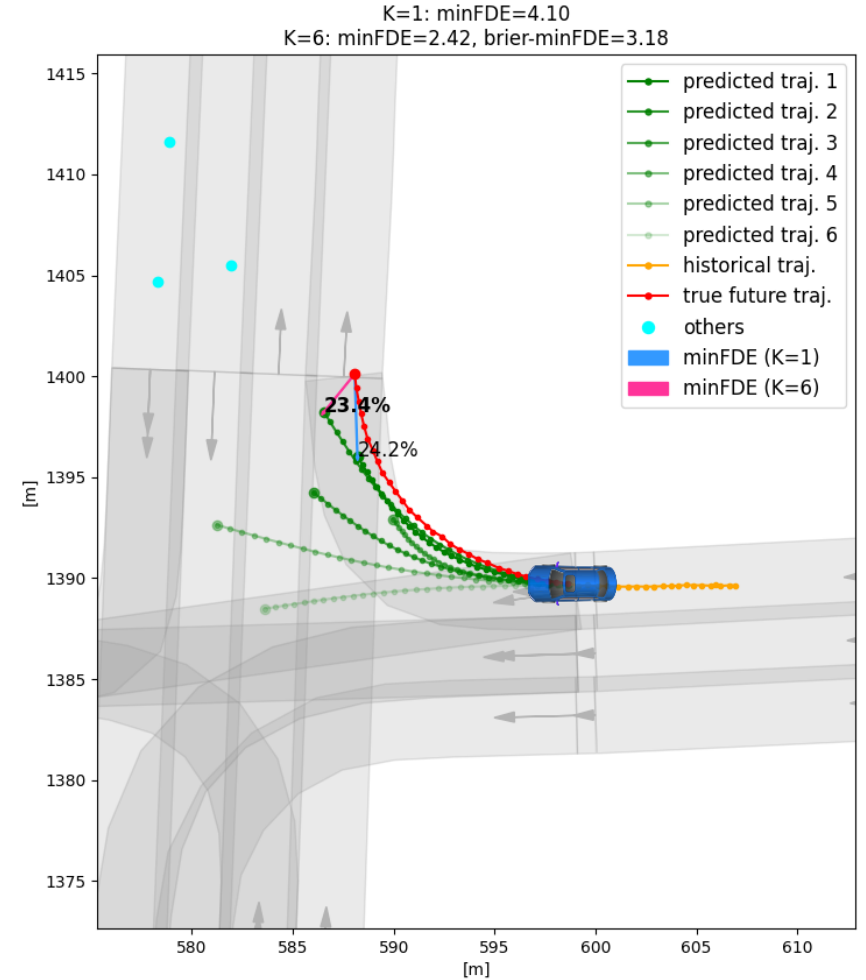
Fig. 26: Turning left

## 3.2 Visualization Results

### Comparison with the original HiVT-64 model



(a) Proposed

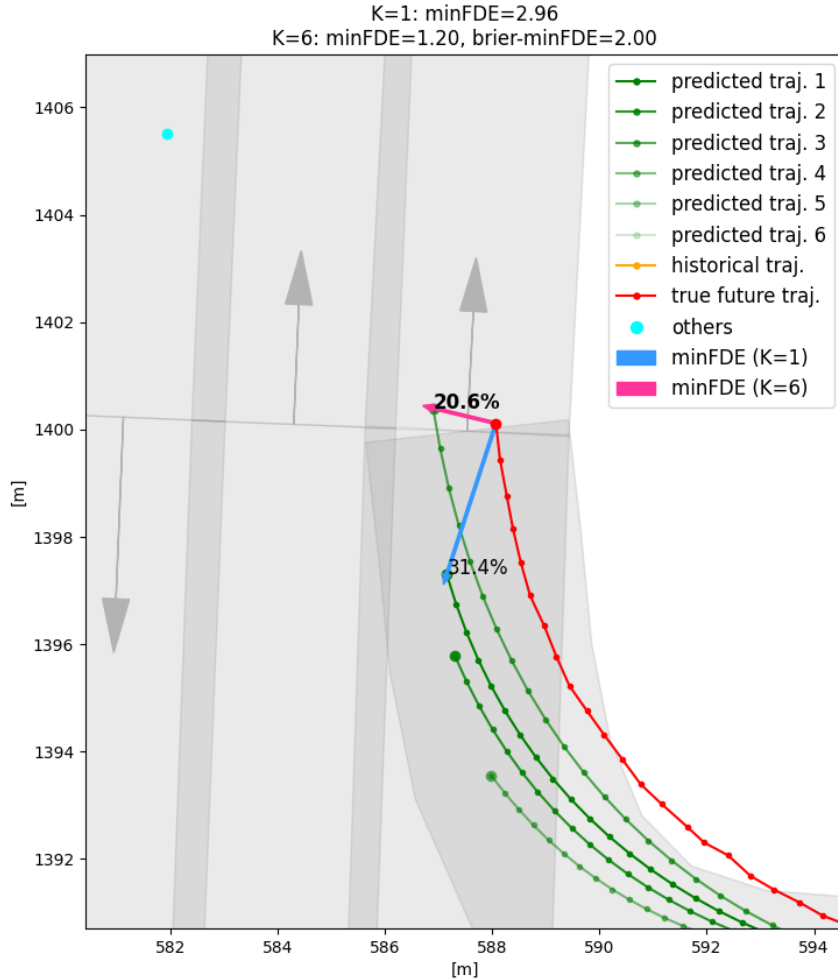


(b) Original

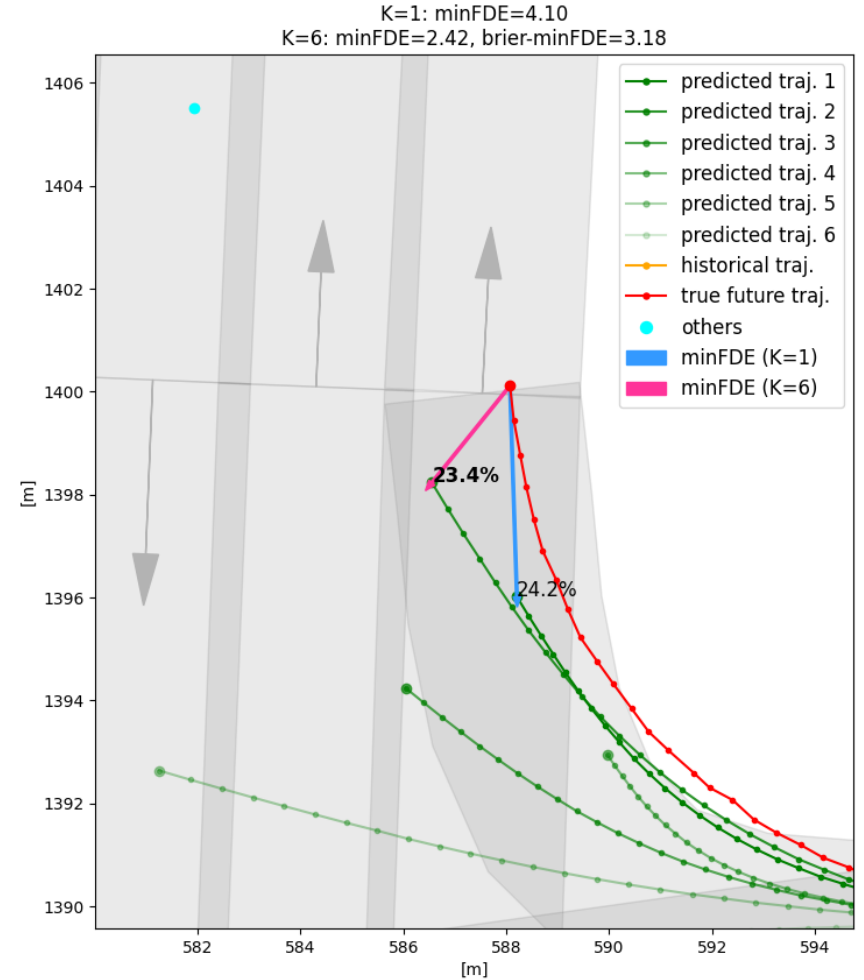
Fig. 27: Turning right (long historical trajectory)

## 3.2 Visualization Results

### Comparison with the original HiVT-64 model (zoom in)



(a) Proposed

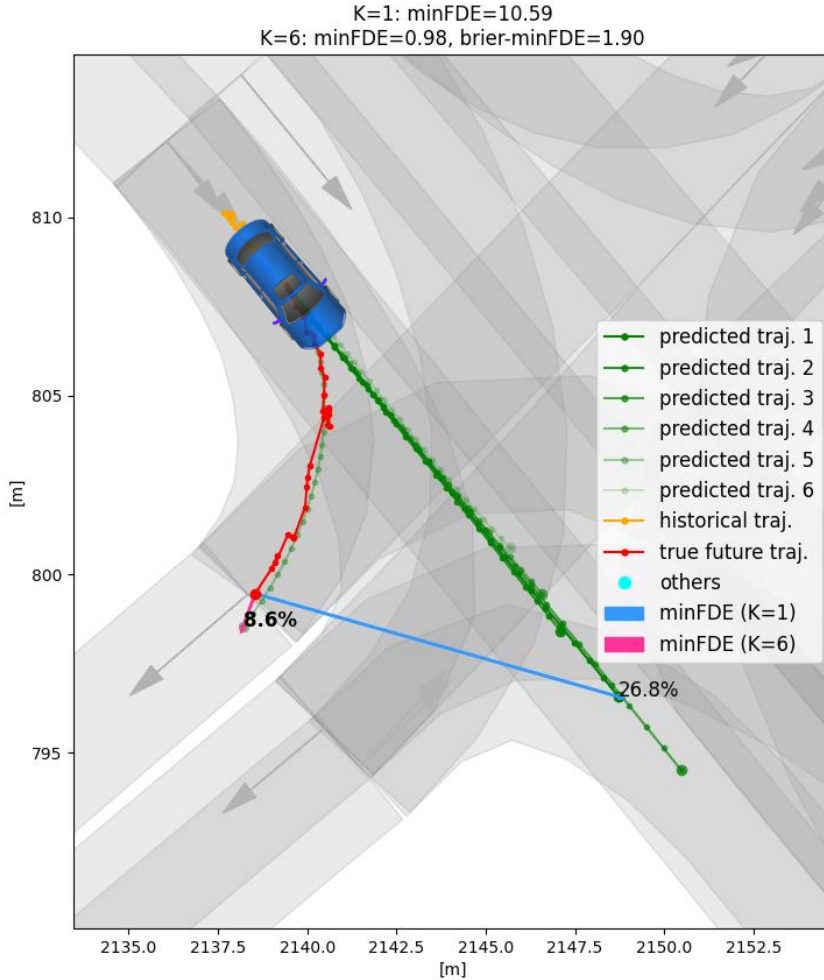


(b) Original

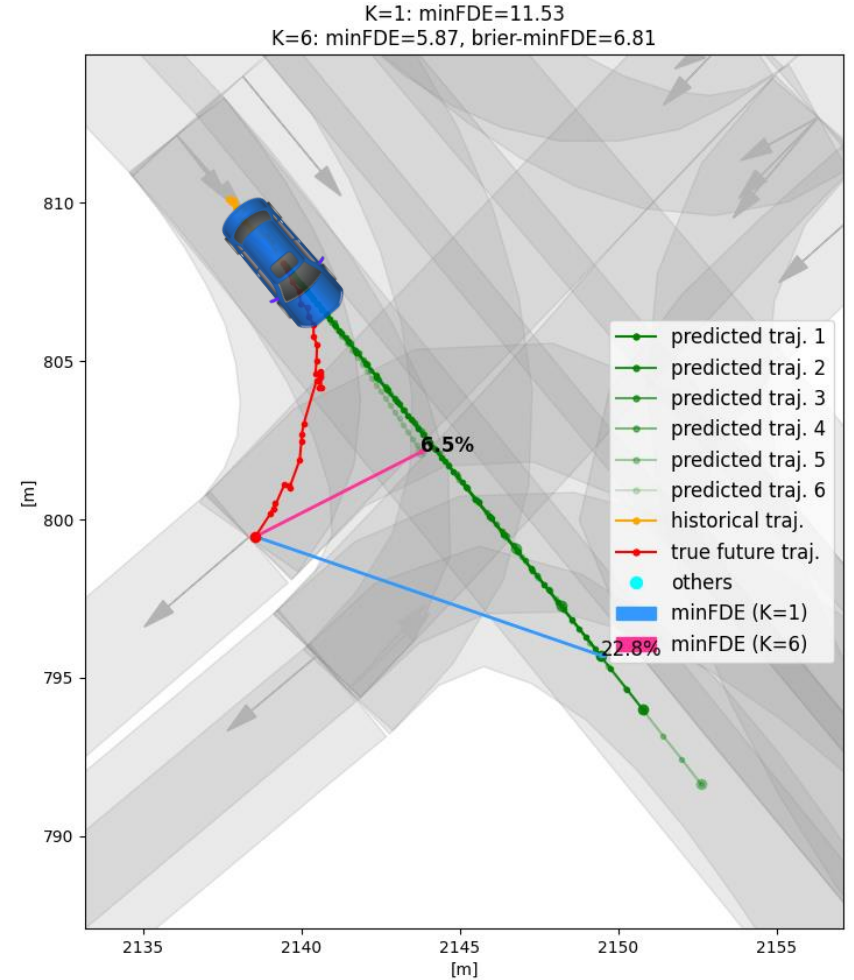
Fig. 27: Turning right (long historical trajectory)

## 3.2 Visualization Results

### Comparison with the original HiVT-64 model



(a) Proposed



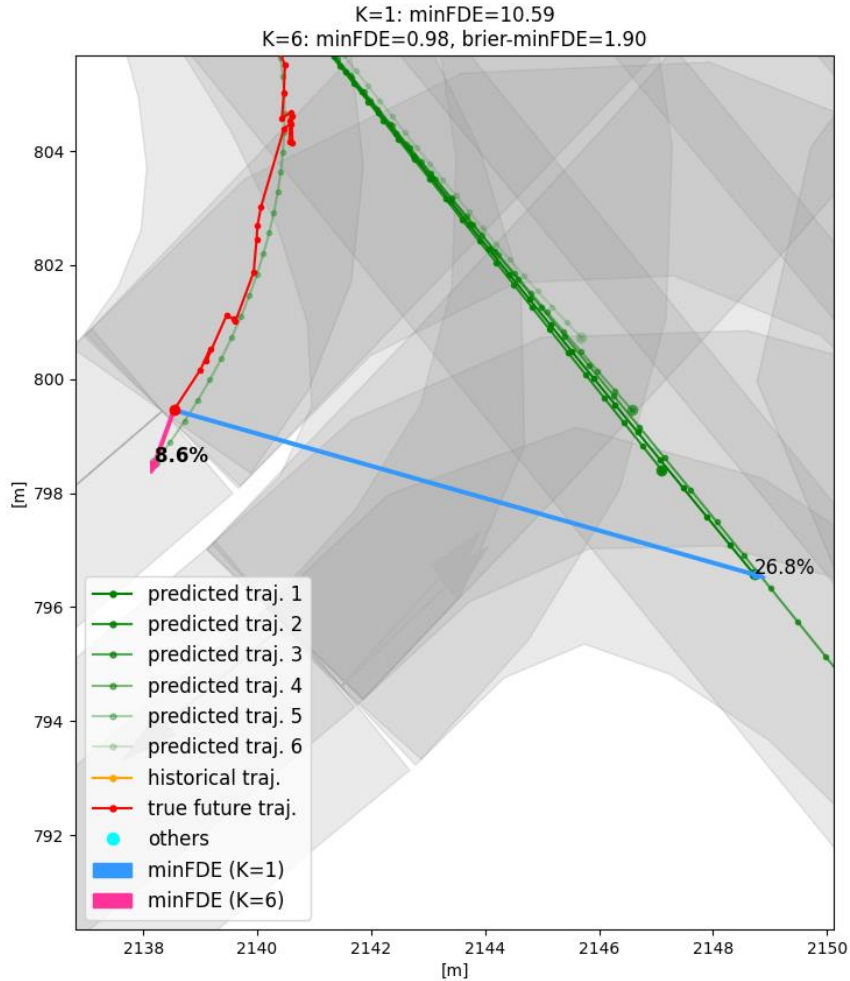
(b) Original

Fig. 28: Turning right (short historical trajectory)

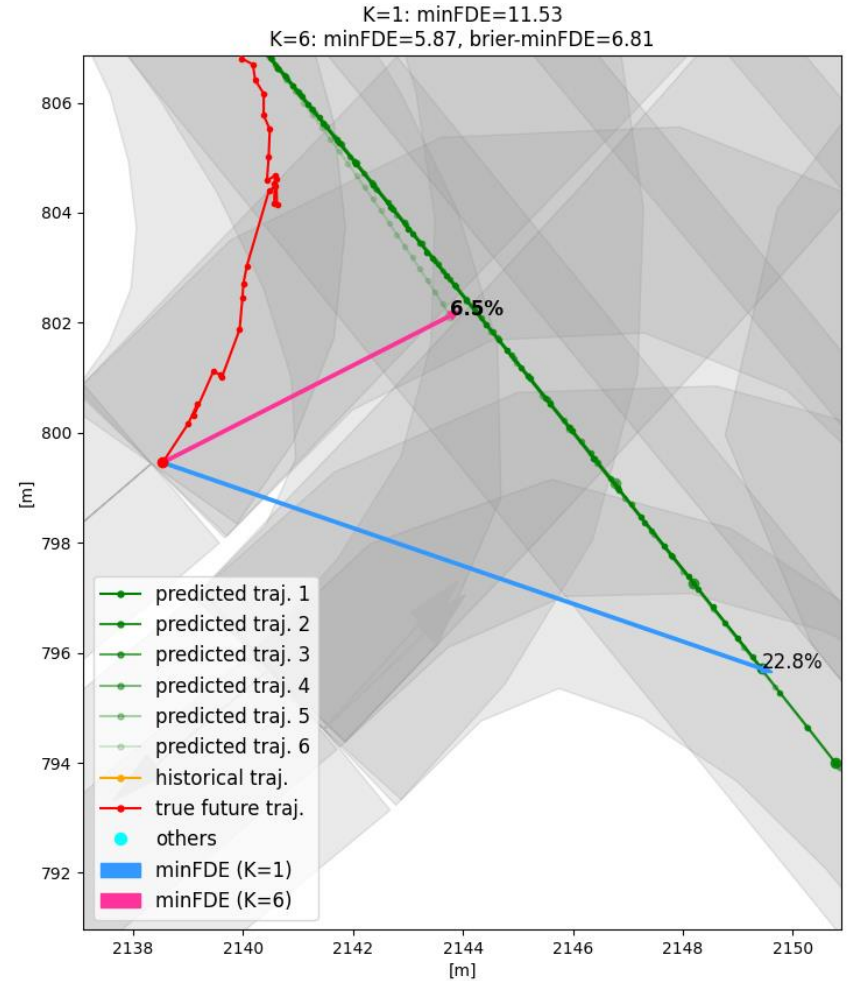


## 3.2 Visualization Results

### Comparison with the original HiVT-64 model (zoom in)



(a) Proposed



(b) Original

Fig. 28: Turning right (short historical trajectory)



# 3.3 Quantitative Comparisons

Metrics Comparison

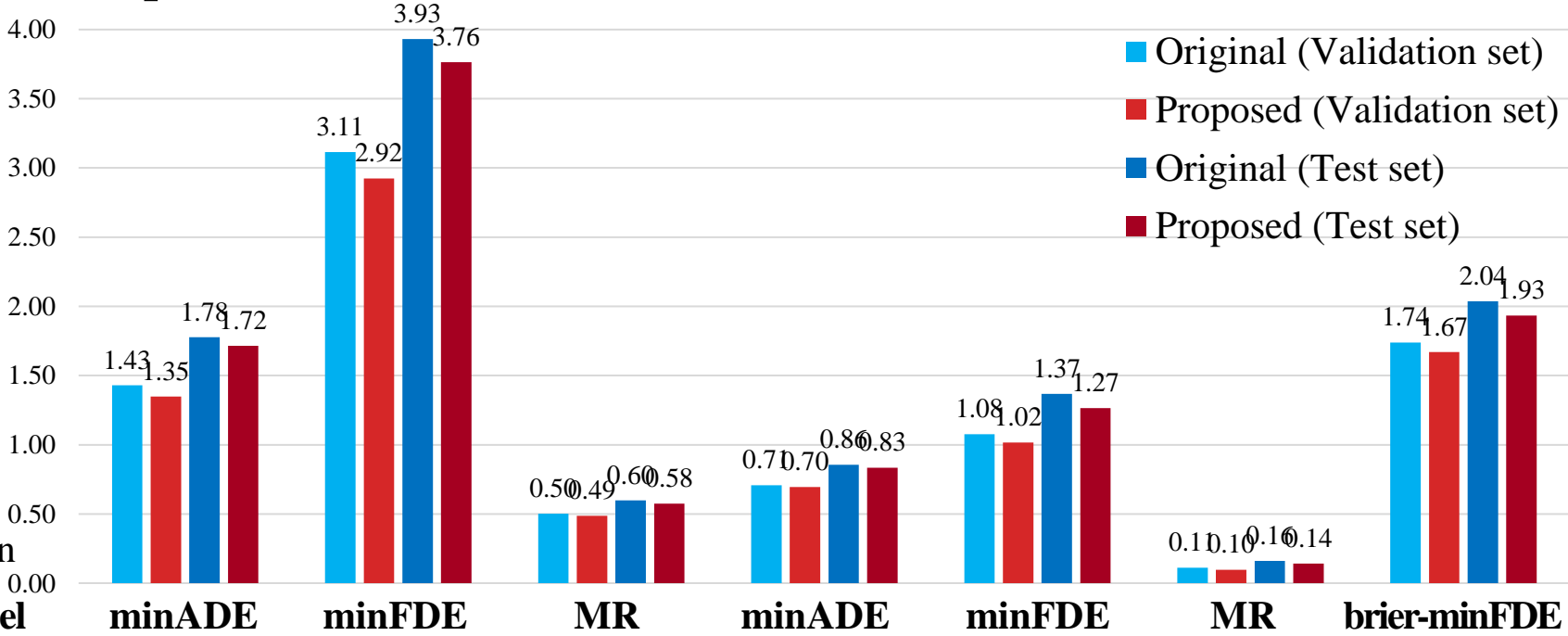


TABLE II:  
Metrics  
comparison

Model		minADE	minFDE	MR	minADE	minFDE	MR	brier-minFDE
Val. Set	Original	1.43	3.11	50.2%	0.71	1.08	11.2%	1.74
	Proposed	1.35	2.92	48.8%	0.70	1.02	9.8%	1.67
	Improve	5.7%	6.1%	2.9%	1.8%	5.6%	11.9%	3.9%
Test set	Original	1.78	3.93	59.9%	0.86	1.37	16.3%	2.04
	Proposed	1.72	3.76	57.6%	0.83	1.27	14.2%	1.93
	Improve	3.4%	4.3%	3.8%	2.4%	7.4%	12.5%	5.1%
#Trajs.		K=1			K=6			

## 3.3 Quantitative Comparisons

### Specification Comparison

#### □ Inference speed (frames per second)

- Test 3 times with the **same 5,000 samples** and take the average

#### □ GPU memory cost

- Feed the **same data** to the network

- Measure the memory cost in **each module**

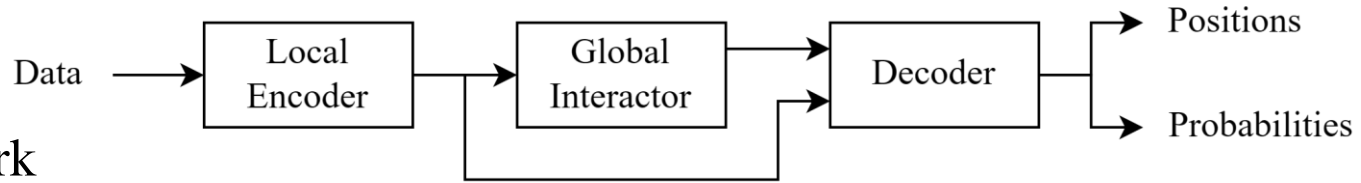


Fig. 30: Rough structure of the model

TABLE III: Specification comparison

Model	#Param.	Inference Speed (fps)	GPU Memory Cost (MB)			
			Local Encoder	Global Interactor	Decoder	Others <sup>†</sup>
Original	653,369	19.11	430	22	4	10
<b>Proposed</b>	<b>395,809</b>	<b>70.70</b>	<b>242</b>	<b>2</b>	<b>4</b>	<b>6</b>
<i>Improve</i>	<i>39.4%</i>	<i>270.0%</i>	<i>43.7%</i>	<i>90.9%</i>	<i>0.0%</i>	<i>40.0%</i>
Lightest*	389,689	74.05	/	/	/	/

\* The model after first-stage modifications.

<sup>†</sup> Includes model and pre-trained weights.

## 3.4 Effectiveness Validation

### Ablation study

TABLE IV: Ablation study of the second-stage modifications (Validation set, K=6)

Decoder Mods.	Encoder Mods.	Training Strategy	minADE	minFDE	MR	brier-minFDE	#Param.
			0.71	1.06	10.9%	1.72	389,689
✓			0.70	1.04	10.3%	1.69	391,117
✓	✓		0.70	1.03	10.1%	1.69	395,809
✓	✓	✓	<b>0.70</b>	<b>1.02</b>	<b>9.8%</b>	<b>1.67</b>	395,809

### Validation of time weights

```
self.time_wgt.squeeze().detach().cpu().numpy()
array([2.062761, 0.51705956, 0.9902031, 1.0632056, 1.1559378,
       1.1959459, 1.2953184, 1.368692, 1.4293575, 1.5335826,
       1.616901, 1.7310858, 1.8082753, 1.9473771, 2.0934815,
       2.3145213, 2.5368783, 2.8813136, 3.2132146, 3.3525586],
      dtype=float32)
```

Fig. 31: Time weights

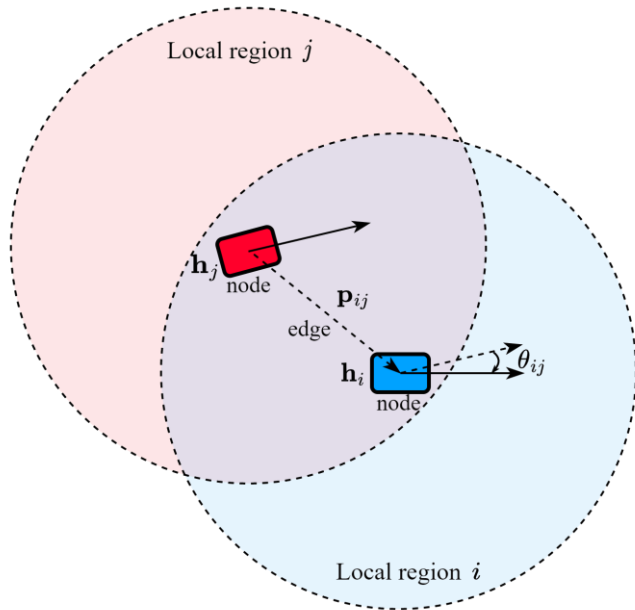


## 3.4 Effectiveness Validation

### Validation of edge attributes fusion

TABLE V: Comparison with the original fusion method (Validation set)

Model	K=1			K=6			
	minADE	minFDE	MR	minADE	minFDE	MR	brier-minFDE
Original	1.37	2.99	49.1%	0.70	1.03	10.1%	1.69
<b>Proposed</b>	<b>1.35</b>	<b>2.92</b>	<b>48.8%</b>	<b>0.70</b>	<b>1.02</b>	<b>9.8%</b>	<b>1.67</b>



#### □ Original

$$\begin{aligned} \blacksquare \mathbf{e}_{ij} &= \phi_{\text{MLP}}([R_{ij}\mathbf{p}_{ij}, \sin(\theta_{ij}), \cos(\theta_{ij})]) \\ \mathbf{h}_j &= \mathbf{h}_j \end{aligned}$$

#### □ Proposed

$$\begin{aligned} \blacksquare \mathbf{e}_{ij} &= [R_{ij}\mathbf{p}_{ij}, \sin(\theta_{ij}), \cos(\theta_{ij})] \\ \mathbf{h}_j &= \phi_{\text{MLP}}([\mathbf{h}_j, \mathbf{e}_{ij}]) \end{aligned}$$

$$\begin{aligned} Q &= W_{Q_{\text{node}}} \mathbf{h}_i \\ K &= W_{K_{\text{node}}} \mathbf{h}_j + W_{K_{\text{edge}}} \mathbf{e}_{ij} \\ V &= W_{V_{\text{node}}} \mathbf{h}_j + W_{V_{\text{edge}}} \mathbf{e}_{ij} \end{aligned}$$

Fig. 17: Illustration of edge attributes



# Contents

---

## 1. Introduction

~~1.1 Background & Research Target~~

1.2 Issue & Idea

---

## 2. Deep Learning-based Trajectory Prediction

2.0 Intro to HiVT (Hierarchical Vector Transformers)

2.1 Computational Efficiency Improvement

2.2 Prediction Accuracy Improvement

---

## 3. Simulations

3.1 Simulation Setup

3.3 Quantitative Comparisons

3.2 Visualization Results

3.4 Effectiveness Validation

---

## 4. Conclusions



# 4. Conclusions

## Conclusions

□ Simulation results: The proposed model **outperforms** the original in both **computational efficiency** and **accuracy**

■ **Computational efficiency** improvement:

□ Inference speed: **2.8x ↑**

Number of parameters: **40.4% ↓**

□ No considerable accuracy degradation

■ **Prediction accuracy** improvement:

□ Inference speed: **2.7x ↑**

Number of parameters: **39.4% ↓**

□ GPU memory cost: **45.5% ↓**

Accuracy\*: **~ 5.5% ↑**

→ The two-stage modifications are **effective**

□ Ablation study: Each modification **contributes** to the **accuracy improvement**

□ Modifications in this research are more **effective on small models**

(Refer to Appendix 6)

■ As the number of parameters increases, the advantage becomes smaller:

□ Hidden size = 64, Accuracy: **~ 5.5% ↑**

□ Hidden size = 128, Accuracy: **~ 4.0% ↑**



**WASEDA UNIVERSITY**

\* Average of all improved percentages



Thanks for Listening



# Appx. 4 Detailed Comparison of Parameters

## Original model

Number of params: 653,369

## Proposed model

Number of params: 395,809

name	#params	name	#params
-----	-----	-----	-----
model	0.7M	model	0.4M
local_encoder	0.4M	local_encoder	0.3M
local_encoder.aa_encoder	86.0K	local_encoder.agt_embed	8.8K
local_encoder.temporal_encoder	0.2M	local_encoder.bos_token	1.3K
local_encoder.al_encoder	76.3K	local_encoder.temporal_encoder	0.1M
global_interactor	0.3M	local_encoder.al_encoder	63.8K
global_interactor.rel_embed	13.4K	local_encoder.aa_encoder	62.7K
global_interactor.global_interactor_layers	0.2M	global_interactor	92.2K
global_interactor.norm	0.1K	global_interactor.global_interactor_layer	67.1K
global_interactor.multihead_proj	25.0K	global_interactor.multihead_proj	25.1K
decoder	37.5K	decoder	30.4K
decoder.aggr_embed	8.4K	decoder.aggr_embed	8.4K
decoder.loc	8.2K	decoder.loc	8.8K
decoder.scale	8.2K	decoder.scale	4.6K
decoder.pi	12.7K	decoder.pi	8.5K



# Appx. 5 Metrics Comparison

## Metrics Comparison

TABLE A5-1: Metrics comparison (Validation set)

Model	minADE	minFDE	MR	minADE	minFDE	MR	brier-minFDE
Original	1.4293	3.1138	0.5022	0.7082	1.0757	0.1117	1.7384
Proposed	<b>1.3480</b>	<b>2.9234</b>	<b>0.4875</b>	<b>0.6955</b>	<b>1.0157</b>	<b>0.0984</b>	<b>1.6707</b>
Improve	5.69%	6.11%	2.92%	1.80%	5.57%	11.91%	3.89%
#Trajs.	K=1			K=6			

TABLE A5-2: Metrics comparison (Test set)

Model	minADE	minFDE	MR	minADE	minFDE	MR	brier-minFDE
Original	1.7771	3.9306	0.5991	0.8552	1.3672	0.1625	2.0372
Proposed	<b>1.7161</b>	<b>3.7626</b>	<b>0.5762</b>	<b>0.8349</b>	<b>1.2654</b>	<b>0.1422</b>	<b>1.9338</b>
Improve	3.43%	4.28%	3.82%	2.38%	7.44%	12.49%	5.07%
#Trajs.	K=1			K=6			



# Appx. 6 Comparison with Other Methods

## Comparison with other methods on the leaderboard [8]

TABLE A6: Comparison with other methods on the leaderboard (Test set, K=6)

Model	minADE	minFDE	MR	brier-minFDE	#Params.
LaneGCN	0.8679	1.3640	0.1634	2.0585	3,701K
DenseTNT	0.8817	1.2815	0.1258	1.9759	1,103K
GOHOME	0.9425	1.4503	0.1048	1.9834	400K
HOME + GOHOME	0.8904	1.2919	0.0846	1.8601	5,100K
Scene Transformer	0.8026	1.2321	0.1255	1.8868	15,296K
HiVT-64	0.8552	1.3672	0.1625	2.0372	653K
HiVT-128	0.8209	1.2800	0.1482	1.9493	2,560K
Proposed	0.8349	1.2654	0.1422	1.9338	396K
Proposed-128	0.8140	1.2213	0.1326	1.8909	1,553K

Hidden size = 64

Hidden size = 128



WASEDA UNIVERSITY

# Appx. 7 Detailed Comparisons

## Ablation study

TABLE A7-1: Ablation study of the second-stage modifications (Validation set, K=6)

Decoder Mods.	Encoder Mods.	Training Strategy	minADE	minFDE	MR	brier-minFDE	#Param.
			0.7110	1.0608	0.1091	1.7236	389,689
✓			0.7042	1.0346	0.1027	1.6935	391,117
✓	✓		0.7041	1.0295	0.1005	1.6852	395,809
✓	✓	✓	<b>0.6955</b>	<b>1.0157</b>	<b>0.0984</b>	<b>1.6707</b>	395,809

## Validation of edge attributes fusion

TABLE A7-2: Comparison with the original fusion method (Validation set)

Model	K=1			K=6			
	minADE	minFDE	MR	minADE	minFDE	MR	brier-minFDE
Original	1.3732	2.9864	0.4913	0.7018	1.0280	0.1010	1.6868
Proposed	<b>1.3480</b>	<b>2.9234</b>	<b>0.4875</b>	<b>0.6955</b>	<b>1.0157</b>	<b>0.0984</b>	<b>1.6707</b>