# A De-Stationary and Cross-Attention LSTM Model for Highway Vehicle Trajectory Prediction

**TianQi Yang[1], Ye Lin[1,2,*]**

[1]College of Mechanical Engineering, Tianjin University of Science and Technology, Tianjin 300222, China
[1,2]Automotive Big Data and Intelligent Technology Laboratory, Tianjin University of Science and Technology, Tianjin 300222, China
*Correspondence Author*

**Abstract:** *Accurate vehicle trajectory prediction is important for autonomous driving, especially in highway scenarios with high-speed motion and complex vehicle interactions. To address these challenges, this paper proposes a de-stationary and cross-attention long short-term memory (DCA-LSTM) model for highway vehicle trajectory prediction. The model combines de-stationary temporal feature enhancement, multi-scale temporal convolution, and ego-centric spatial cross-attention to improve temporal modeling and interaction representation. The de-stationary temporal module is designed to improve the representation of motion sequences with changing statistical properties, while the multi-scale convolution structure helps capture both short-term fluctuations and long-term motion trends. In addition, the spatial cross-attention mechanism dynamically aggregates interaction information from neighboring vehicles according to their influence on the target vehicle. Experiments on the NGSIM US-101 and I-80 datasets show that the proposed method outperforms the compared baseline models across multiple prediction horizons and achieves lower RMSE, especially for 5 s prediction. The ablation and qualitative results further demonstrate the effectiveness of the proposed temporal enhancement and spatial interaction modeling modules. The results indicate that the proposed model is effective for long-horizon highway trajectory prediction in complex traffic environments.*

**Keywords:** Autonomous driving, Vehicle trajectory prediction, LSTM, Cross-attention, Highway driving.

## 1. Introduction

Vehicle trajectory prediction is an important component of autonomous driving because it provides essential support for decision-making and motion planning [1-2]. An autonomous vehicle operating in dynamic traffic environments must not only perceive the current states of surrounding participants, but also estimate their future movements in advance. This capability allows the vehicle to identify potential conflicts, adjust driving strategies, and generate safer and smoother control actions. In highway scenarios, the role of trajectory prediction becomes even more critical because vehicles usually travel at relatively high speeds, lane-level interactions may evolve quickly, and safety margins are often limited.

Highway driving environments present several characteristics that make trajectory prediction difficult. Compared with low-speed urban environments, highway traffic usually exhibits stronger continuity in motion while still containing sudden changes caused by overtaking, car-following adjustments, lane changes, and surrounding traffic disturbances. As a result, a prediction model must be able to describe both stable motion tendencies and abrupt local variations. Moreover, the future motion of a target vehicle is not determined solely by its own historical trajectory. It is also influenced by neighboring vehicles in front, behind, and in adjacent lanes. The relative importance of these neighboring vehicles may change over time, which further increases the difficulty of accurate prediction.

Another major challenge is the non-stationary nature of highway motion. In practice, the statistical properties of vehicle motion sequences may vary across time because of changes in speed, acceleration, driver intention, and traffic context [3–6]. A model that assumes relatively stable temporal patterns may perform well for short prediction horizons but become less reliable for longer horizons, where motion evolution becomes more uncertain. In addition, highway

trajectories usually contain motion information at multiple temporal scales. Long-term trends reflect maneuver evolution and driving objectives, whereas short-term variations capture local motion adjustments and immediate responses to surrounding traffic. Therefore, effective trajectory prediction requires a model that can simultaneously account for temporal non-stationarity, interaction heterogeneity, and multi-scale motion characteristics.

Existing trajectory prediction methods can generally be divided into three categories: physics-based methods, traditional machine learning methods, and deep learning methods. Physics-based methods rely on simplified motion assumptions, such as constant velocity or constant acceleration, and therefore are computationally efficient and easy to implement [7-8]. However, these methods often have limited performance in long-horizon prediction because they cannot adequately describe complex maneuvers or interaction-aware motion. Traditional machine learning methods can learn motion patterns from historical data and have shown advantages over purely model-based approaches in some scenarios [9], [10]. Nevertheless, such methods usually depend on handcrafted features and often have limited capability in modeling high-dimensional spatio-temporal dependencies.

Deep learning methods, especially those based on recurrent neural networks and Long Short-Term Memory (LSTM) networks, have shown good performance in trajectory prediction tasks because they are able to model sequential dependencies in historical trajectories [11]. To further improve prediction accuracy, researchers have introduced convolutional structures, social pooling mechanisms, and attention mechanisms to capture local patterns and vehicle interactions [12–15]. Although these methods have improved trajectory prediction to some extent, some limitations remain.

In particular, many existing methods do not explicitly address

the non-stationary nature of highway trajectories, and some interaction modeling strategies are not sufficiently flexible in distinguishing the different influence levels of surrounding vehicles. In addition, temporal features at different granularities are often not fully captured within a unified framework.

To address these issues, this paper proposes a highway vehicle trajectory prediction model termed DCA-LSTM, which combines de-stationary temporal enhancement with spatial cross-attention. The proposed model introduces a de-stationary temporal feature enhancement mechanism and a multi-scale temporal convolution module to improve temporal representation under changing motion statistics. Meanwhile, an ego-centric spatial cross-attention mechanism is employed to model target-neighbor interactions more effectively by assigning adaptive importance to surrounding vehicles. Experiments on the NGSIM US-101 and I-80 datasets show that the proposed method achieves better prediction performance than the compared baseline methods across multiple prediction horizons.

The main contributions of this paper can be summarized as follows. First, a de-stationary temporal enhancement mechanism is introduced to improve the modeling of trajectory sequences with changing statistical properties. Second, multi-scale temporal convolution is incorporated to capture motion patterns at different temporal granularities. Third, an ego-centric spatial cross-attention mechanism is designed to represent the influence of neighboring vehicles more effectively. Fourth, extensive experiments on real highway trajectory datasets verify the effectiveness of the proposed method in quantitative, ablation, and qualitative evaluations.

The remainder of this paper is organized as follows. Section 2 presents the proposed methodology. Section 3 describes the experiments and discusses the results. Section 4 concludes the paper.

## 2. Methodology

### 2.1 Problem Formulation

Vehicle trajectory prediction aims to estimate the future positions of a target vehicle based on its observed historical motion and surrounding traffic information. Let the state of vehicle i at time t be denoted as $x_i^t$. In a highway traffic scene, the state information may include position-related variables and motion-related variables, such as lateral position, longitudinal position, velocity, and acceleration. Given the historical trajectory sequence X over an observation horizon $T_{obs}$, the objective is to learn a mapping function f that predicts the future trajectory of the target vehicle over a prediction horizon $T_{pred}$:

$$\hat{Y} = f(X) = \left\{ \widehat{p_{ego}^{t+1}}, \ldots, \widehat{p_{ego}^{t+T_{pred}}} \right\} \qquad (1)$$

where $\widehat{p_{ego}^t}$ denotes the predicted position of the target vehicle at time t.

In this study, the model input includes the historical motion information of the target vehicle and surrounding vehicles, while the output is the future trajectory of the target vehicle. This formulation reflects the practical requirement of highway autonomous driving, where the future behavior of the target vehicle is closely influenced by both its own historical motion and the surrounding traffic context. The prediction task is therefore essentially a spatio-temporal sequence modeling problem.

### 2.2 Overall Architecture of DCA-LSTM

The proposed DCA-LSTM adopts an encoder-decoder architecture. It consists of four main components: feature embedding, temporal encoding with de-stationary enhancement, spatial interaction modeling, and trajectory decoding. As shown in Figure 1, the framework is designed to improve temporal feature representation and interaction modeling in highway scenarios.
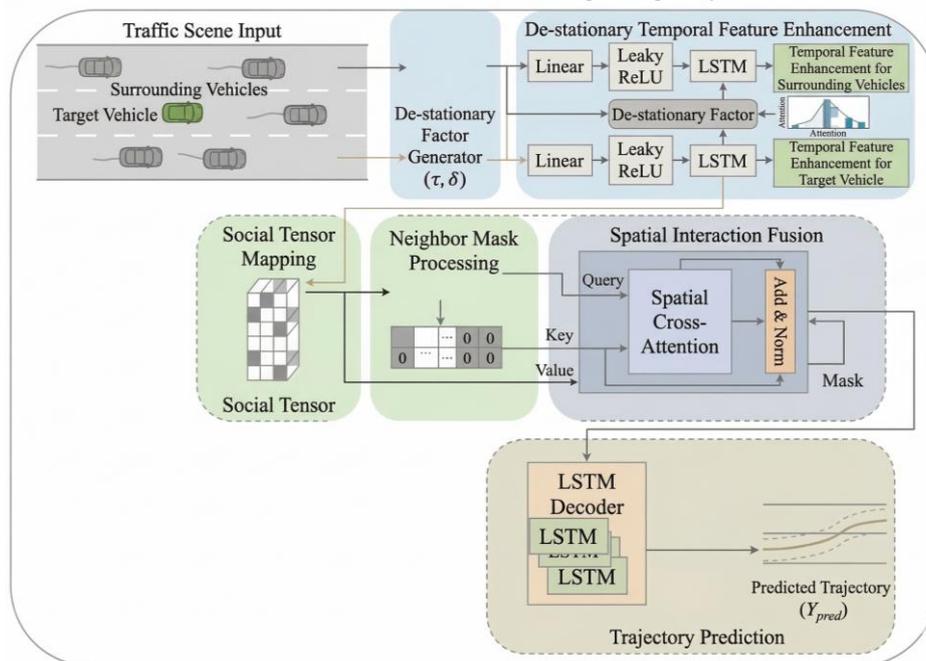


**Figure 1:** Overall architecture of DCA-LSTM

The overall processing pipeline can be summarized as follows. First, the historical trajectory features of the target vehicle and surrounding vehicles are embedded into a latent representation space. This step helps the network transform raw motion variables into features that are more suitable for subsequent temporal and spatial modeling. Second, the temporal branch extracts motion-related information from the historical sequences. In this stage, the de-stationary enhancement mechanism and the multi-scale temporal convolution module are used to strengthen temporal representation. Third, spatial interaction information is incorporated through an ego-centric spatial cross-attention mechanism, in which the target vehicle interacts with surrounding vehicles in an adaptive manner. Finally, the fused spatio-temporal representation is fed into the decoder to generate future trajectory predictions.

Compared with a conventional LSTM-based prediction framework, the proposed architecture emphasizes two aspects. On the one hand, it explicitly considers the changing statistical properties of motion sequences. On the other hand, it models the influence of surrounding vehicles through adaptive interaction weighting rather than simple aggregation. In this way, the model jointly considers temporal evolution and surrounding vehicle influence within a unified framework.

## 2.3 De-stationary Temporal Feature Enhancement

Highway vehicle trajectories often exhibit non-stationary characteristics because the motion state of a vehicle may change under different driving conditions. For example, a vehicle may maintain a relatively stable speed during lane keeping, while sudden changes in speed or acceleration may occur during car-following adjustment or overtaking. These changes can cause the statistical properties of trajectory sequences to vary over time, which may reduce the effectiveness of conventional temporal attention mechanisms.

To address this issue, a de-stationary temporal enhancement mechanism is introduced into the temporal attention module. Specifically, adjustment factors $\tau$ and $\delta$ are used to refine the attention scores:

$$\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\tau\sqrt{d_k}} + \delta\right)V \qquad (2)$$

where $\tau$ and $\delta$ are used to adjust the scale and shift of the attention distribution. This design helps reduce the influence of changing motion statistics on temporal representation and allows the attention mechanism to better adapt to different motion states.

In addition, a multi-scale temporal convolution module with parallel kernels of different sizes is employed to capture motion patterns at different temporal granularities. Small kernels are helpful for describing local motion changes and short-term fluctuations, while larger kernels are more suitable for modeling broader motion trends over longer time spans. By combining multiple temporal scales, the model can better represent both short-term and long-term information in the historical trajectory sequence.

As shown in Figure 2, this module enhances temporal representation by combining de-stationary adjustment and multi-scale feature extraction. The purpose of this design is to improve the robustness of temporal modeling when vehicle motion varies across different driving states. This is particularly important for long-horizon prediction, where small modeling errors in temporal evolution may accumulate over time
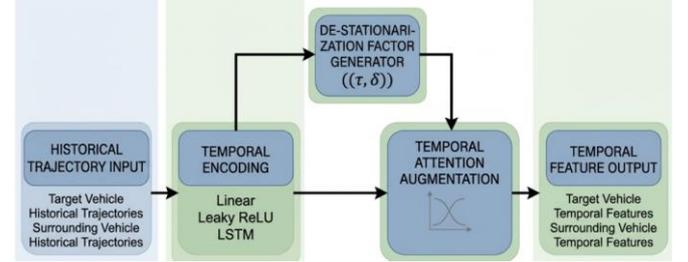


**Figure 2:** De-stationary temporal feature enhancement module

## 2.4 Spatial Interaction Modeling

In highway traffic, the future motion of a target vehicle is strongly influenced by surrounding vehicles. For example, a leading vehicle may determine car-following behavior, while a neighboring vehicle in an adjacent lane may affect lateral motion decisions. Therefore, accurate trajectory prediction requires a model that can represent these interactions effectively.

To model interactions between the target vehicle and surrounding vehicles, an ego-centric spatial cross-attention mechanism is employed. In this module, the target vehicle feature is used as the query, while neighboring vehicle features are used as keys and values. The attention weight for each neighboring vehicle is computed as:

$$\alpha_{i,j} = \frac{\exp\left(\text{score}(Q_i, K_j)\right)}{\sum_{k \in N_i} \exp\left(\text{score}(Q_i, K_k)\right)} \qquad (3)$$

where $N_i$ denotes the neighboring vehicle set of the target vehicle. This mechanism enables the model to focus on surrounding vehicles that have stronger influence on the future motion of the target vehicle.

The ego-centric design is particularly suitable for trajectory prediction because the main objective is to predict the motion of the target vehicle rather than to represent the scene in a uniform manner. By using the target vehicle as the interaction center, the model can selectively aggregate information from surrounding vehicles according to their relevance. Compared with simple feature aggregation or grid-based pooling, the proposed mechanism can assign different importance to neighboring vehicles according to their effect on the target vehicle.

As shown in Figure 3, the proposed spatial interaction module captures the influence of neighboring vehicles in an ego-centric manner. This design helps the network better represent dynamic interaction relationships in highway environments and improves the ability to model interaction-dependent trajectory evolution.
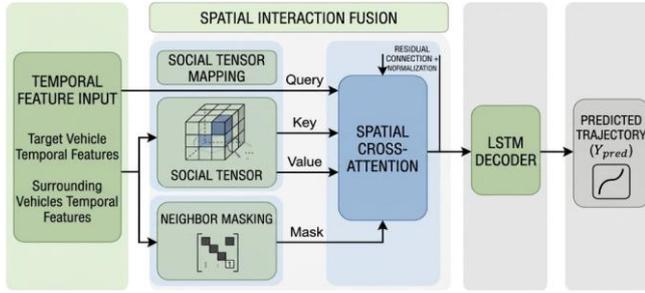
**Figure 3:** Spatial cross-attention interaction modeling

## 2.5 Decoder and Loss Function

After temporal and spatial features are fused, the resulting spatio-temporal representation is fed into an LSTM decoder to generate the predicted future trajectory. The decoder transforms the learned representation into future positions over the prediction horizon. In this way, the model connects historical spatio-temporal information with future motion estimation in a sequential prediction framework.

During training, the model is optimized using a masked mean squared error (MSE) loss:

$$\mathcal{L} = \frac{1}{M} \sum_{i=1}^{N} \sum_{t=1}^{T_{pred}} l_{i,t} \cdot |p_{i,t} - \widehat{p_{i,t}}|^2 \qquad (4)$$

where $l_{i,t}$ is the mask for valid trajectory points, and M denotes the number of valid samples. By minimizing this loss, the model is trained to reduce the deviation between predicted trajectories and ground-truth trajectories over the prediction horizon.

The use of a masked loss helps ensure that only valid trajectory points contribute to training. This is important for trajectory datasets in which some samples may have missing or invalid points in certain time steps. Overall, the decoder and loss function provide a straightforward training objective while allowing the proposed temporal and spatial modules to contribute directly to prediction accuracy.

## 3. Experiments and Discussion

### 3.1 Dataset and Preprocessing

The proposed model is evaluated on the NGSIM dataset, including the US-101 and I-80 highway segments. The dataset is sampled at 10 Hz and contains vehicle position, velocity, and acceleration information. As shown in Figure 4, the dataset covers typical highway traffic scenarios with multi-vehicle interactions.

The NGSIM dataset has been widely used in trajectory prediction studies because it provides detailed and high-frequency highway traffic trajectories. It includes a variety of driving situations such as car-following, lane keeping, and lane-level interactions, which makes it suitable for evaluating trajectory prediction methods in realistic multi-vehicle scenarios. In particular, the presence of dense traffic and frequent local interactions makes this dataset useful for verifying whether a prediction model can handle both temporal evolution and spatial interaction complexity.

To reduce noise in the raw trajectory data, a Savitzky-Golay (SG) smoothing filter is applied with a window length of 11 and a polynomial degree of 3. This preprocessing step suppresses high-frequency noise while preserving the main motion trend. Such smoothing is necessary because trajectory data extracted from video or sensor sources may contain local fluctuations that do not reflect actual vehicle behavior. After preprocessing, the trajectory sequences are more suitable for model training and evaluation.
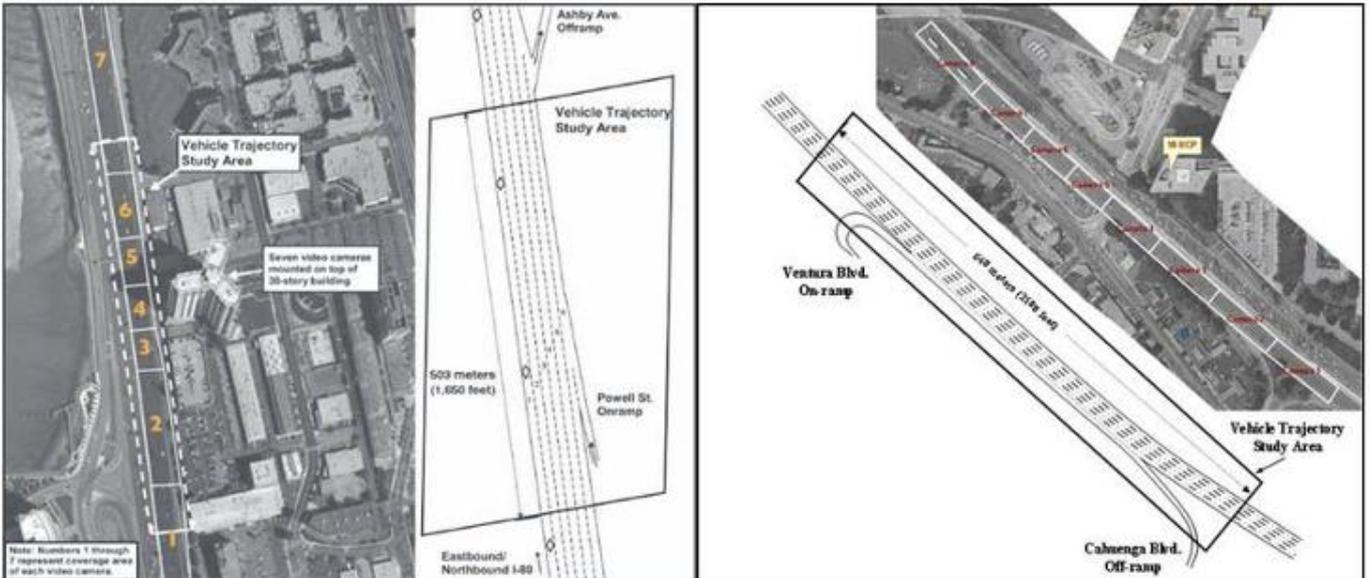


**Figure 4:** Highway traffic scenarios in the NGSIM dataset

### 3.2 Evaluation Metric and Baselines

Root Mean Square Error (RMSE) is used as the evaluation metric. It measures the deviation between the predicted trajectory and the ground truth at each prediction horizon:

$$\text{RMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^{N} \left[ \left( \widehat{x_i^t} - x_i^t \right)^2 + \left( \widehat{y_i^t} - y_i^t \right)^2 \right]} \qquad (5)$$

where N denotes the number of test samples. RMSE is a widely used metric in trajectory prediction because it directly

reflects the average positional error between predicted trajectories and real ground-truth trajectories.

The proposed method is compared with several baseline models, including standard LSTM, CNN-LSTM, Convolutional Social LSTM (CS-LSTM), and Spatio-Temporal Attention LSTM (STA-LSTM). These baselines are selected because they represent different levels of temporal modeling and interaction modeling capability. The standard LSTM model serves as a basic sequential prediction baseline. CNN-LSTM introduces convolutional operations to strengthen local feature extraction. CS-LSTM emphasizes social interaction modeling through structured feature aggregation, while STA-LSTM incorporates attention mechanisms to model spatio-temporal relationships more adaptively. Therefore, these baselines provide a meaningful comparison framework for evaluating the proposed DCA-LSTM.

### 3.3 Quantitative Results and Discussion

As shown in Table 1, the proposed DCA-LSTM achieves the lowest RMSE across different prediction horizons.

**Table 1:** RMSE comparison of different models under different prediction horizons

| Evaluation Metric | Prediction Horizon | LSTM | CNN-LSTM | CS-LSTM | STA-LSTM | DCA-LSTM |
|---|---|---|---|---|---|---|
| | 1s | 0.88 | 0.61 | 0.61 | 0.58 | 0.45 |
| | 2s | 1.54 | 1.34 | 1.27 | 1.23 | 1.05 |
| RMSE (m) | 3s | 2.55 | 2.27 | 2.09 | 2.01 | 1.78 |
| | 4s | 4.02 | 3.42 | 3.10 | 3.04 | 2.68 |
| | 5s | 5.96 | 4.95 | 4.36 | 4.27 | 3.78 |

The results show that models incorporating interaction information, including CS-LSTM, STA-LSTM, and DCA-LSTM, generally outperform models without explicit spatial interaction modeling. Among them, DCA-LSTM achieves the best performance at all prediction horizons. This indicates that the proposed method is more effective in modeling temporal motion patterns and vehicle interactions.

For short-term prediction, most models achieve relatively small errors because recent motion information is still highly informative. In such cases, the future trajectory is strongly constrained by the immediate historical motion of the target vehicle, and even relatively simple temporal models can provide acceptable predictions. However, as the prediction horizon increases, the error gap between different models becomes more apparent. This is because longer prediction horizons require stronger capability in preserving motion trends and handling trajectory evolution under changing traffic conditions.

The advantage of DCA-LSTM becomes more evident in long-horizon prediction. This suggests that the proposed de-stationary temporal enhancement and multi-scale temporal modeling are helpful for maintaining stable temporal representation when the motion statistics vary over time. At the same time, the ego-centric spatial cross-attention mechanism provides more effective interaction modeling, which helps the network capture how surrounding vehicles influence the future motion of the target vehicle. Therefore, the superior long-horizon performance of DCA-LSTM can be understood as the result of jointly improving temporal robustness and interaction sensitivity.

From another perspective, the quantitative results also show that highway trajectory prediction cannot rely solely on target-vehicle historical motion. Models that explicitly account for interaction information consistently perform better, which confirms the importance of spatial context in realistic highway traffic scenarios. The proposed DCA-LSTM further improves this performance by using a more flexible interaction mechanism and a more robust temporal representation strategy.

### 3.4 Ablation Study

To further evaluate the contribution of each module, ablation experiments are conducted on the proposed model. As shown in Table 2, removing either the temporal enhancement module or the spatial interaction module leads to performance degradation.

**Table 2:** Ablation study results under different prediction horizons

| Evaluation Metric | Prediction Horizon | w/o Temporal Attention | w/o Spatial Attention | DCA-LSTM |
|---|---|---|---|---|
| | 1s | 0.53 | 0.52 | 0.45 |
| | 2s | 2.12 | 2.09 | 1.05 |
| RMSE(m) | 3s | 3.04 | 2.99 | 1.78 |
| | 4s | 3.65 | 3.61 | 2.68 |
| | 5s | 4.29 | 4.08 | 3.78 |

The ablation results indicate that both temporal enhancement and spatial interaction modeling contribute to the final prediction performance. Temporal enhancement improves the representation of non-stationary motion, while spatial interaction modeling improves the ability to capture the influence of surrounding vehicles. The full DCA-LSTM model achieves the best results, which verifies the effectiveness of combining de-stationary temporal enhancement with spatial cross-attention.

More specifically, when the temporal enhancement component is removed, the model becomes less capable of handling changing motion statistics across the observation sequence. This leads to larger prediction errors, especially for longer horizons where temporal modeling quality has a stronger impact on the final result. When the spatial interaction module is removed or weakened, the model becomes less effective in distinguishing the influence of neighboring vehicles, which reduces its ability to capture interaction-dependent motion changes.

These ablation findings are consistent with the design objectives of the proposed framework. The temporal module is mainly responsible for improving sequence representation under non-stationary motion conditions, while the spatial module is responsible for identifying relevant neighboring

vehicle influence. Their combination leads to the best overall performance, suggesting that both aspects are necessary for highway trajectory prediction in complex traffic environments.

### 3.5 Qualitative Analysis

To further illustrate the prediction performance, a typical trajectory prediction example is shown in Figure 5, the predicted trajectory of DCA-LSTM is closer to the ground truth than that of STA-LSTM, especially in the later stage of the prediction horizon.
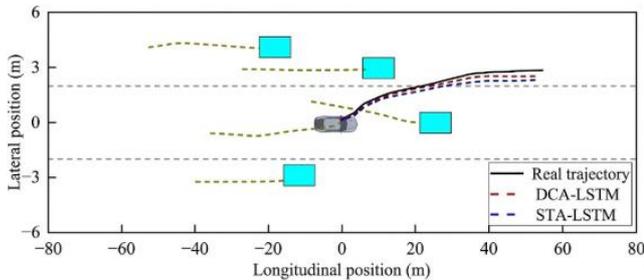


**Figure 5:** Visualization of predicted trajectories

This result shows that the proposed model can better track the motion trend of the target vehicle and capture the interaction influence of surrounding vehicles. In addition, the predicted trajectory of DCA-LSTM shows better consistency with the ground truth in both overall trend and local position changes. The qualitative result is consistent with the quantitative comparison in Table 1.

From the perspective of visual trajectory fitting, the advantage of the proposed method lies not only in reducing numerical error but also in preserving realistic motion evolution. In trajectory prediction tasks, especially in highway scenarios, a good prediction should follow the general motion trend while remaining close to the actual vehicle path. The prediction produced by DCA-LSTM better satisfies these requirements, which suggests that the proposed model can represent both temporal progression and interaction-dependent motion changes more effectively than the compared attention-based baseline.

## 4. Conclusion

This paper proposed a highway vehicle trajectory prediction model termed DCA-LSTM. The model combines de-stationary temporal enhancement, multi-scale temporal convolution, and ego-centric spatial cross-attention to improve temporal representation and interaction modeling. Experimental results on the NGSIM US-101 and I-80 datasets show that the proposed method outperforms the compared baseline models across multiple prediction horizons. The ablation and qualitative results further verify the effectiveness of the proposed modules. Therefore, the proposed method can provide useful support for downstream trajectory planning and decision-making in autonomous highway driving.

## References

[1] S. K. Ahmed et al., "Road traffic accidental injuries and deaths: A neglected global health issue," Health Sci. Rep., vol. 6, no. 5, p. e1240, 2023.

[2] P. Koopman and M. Wagner, "Autonomous vehicle safety: An interdisciplinary challenge," IEEE Intell. Transp. Syst. Mag., vol. 9, no. 1, pp. 90–96, 2017.

[3] V. S. R. Kosuru and A. K. Venkitaraman, "Advancements and challenges in achieving fully autonomous self-driving vehicles," World J. Adv. Res. Rev., vol. 18, no. 1, pp. 161–167, 2023.

[4] S. Teng et al., "Motion planning for autonomous driving: The state of the art and future perspectives," IEEE Trans. Intell. Veh., vol. 8, no. 6, pp. 3692–3711, 2023.

[5] H. Li et al., "A physical law constrained deep learning model for vehicle trajectory prediction," IEEE Internet Things J., vol. 10, no. 24, pp. 22775–22790, 2023.

[6] F. Mo, J. Hao, H. Huang, and X. Li, "Trajectory Prediction for Autonomous Driving based on Kinematic and Intention Awareness," IEEE Transactions on Intelligent Transportation Systems, vol. 24, no. 5, pp. 4887-4901, 2023.

[7] Y. Huang, J. Wang, C. Galbraith, K. Herold, and C. Guo, "A Survey on Trajectory-Prediction Methods for Autonomous Driving," IEEE Transactions on Intelligent Vehicles, vol. 7, no. 3, pp. 652-674, 2022.

[8] S. Mozaffari, O. Y. Al-Jarrah, M. Dianati, P. Jennings, and A. Mouzakitis, "Deep Learning for Trajectory Prediction: A Survey," IEEE Transactions on Intelligent Transportation Systems, vol. 22, no. 5, pp. 2563-2575, 2021.

[9] X. Wang, Y. Zheng, and X. Li, "Vehicle Trajectory Prediction Using Data-Driven Machine Learning Methods in Complex Traffic," IEEE Access, vol. 10, pp. 12345-12356, 2022.

[10] Z. Zhao, et al., "A Hybrid Machine Learning Framework for Vehicle Trajectory Prediction in Highway Scenarios," IEEE Transactions on Intelligent Transportation Systems, vol. 24, no. 8, pp. 8320-8332, 2023.

[11] Z. Sheng, Y. Xu, S. Xue, and D. Li, "Graph-based Spatial-Temporal Convolutional Network for Vehicle Trajectory Prediction," IEEE Intelligent Transportation Systems Magazine, vol. 15, no. 1, pp. 150-162, 2023.

[12] N. Deo and M. M. Trivedi, "Convolutional social pooling for vehicle trajectory prediction," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. Workshops, 2018, pp. 1468–1476.

[13] L. Lin et al., "Vehicle trajectory prediction using LSTMs with spatial-temporal attention mechanisms," IEEE Intell. Transp. Syst. Mag., vol. 14, no. 2, pp. 197–208, 2021.

[14] S. Hochreiter and J. Schmidhuber, "Long short-term memory," Neural Comput., vol. 9, no. 8, pp. 1735–1780, 1997.

[15] A. Savitzky and M. J. E. Golay, "Smoothing and differentiation of data by simplified least squares procedures," Anal. Chem., vol. 36, no. 8, pp. 1627–1639, 1964.