

TransParking: A Dual-Decoder Transformer Framework with Soft Localization for End-to-End Automatic Parking

Hangyu Du, Chee-Meng Chew*

Abstract—In recent years, fully differentiable end-to-end autonomous driving systems have become a research hotspot in the field of intelligent transportation. Among various research directions, automatic parking is particularly critical as it aims to enable precise vehicle parking in complex environments. In this paper, we present a purely vision-based transformer model for end-to-end automatic parking, trained using expert trajectories. Given camera-captured data as input, the proposed model directly outputs future trajectory coordinates. Experimental results demonstrate that the various errors of our model have decreased by approximately 50% in comparison with the current state-of-the-art end-to-end trajectory prediction algorithm of the same type. Our approach thus provides an effective solution for fully differentiable automatic parking.

I. INTRODUCTION

The traditional automatic parking systems are usually composed of three modules: perception, planning and control. In these systems, the vehicle initially employs sensors such as cameras and LiDAR to detect and model its surrounding environment, identifying both the target parking space and any nearby obstacles. Subsequently, a trajectory is planned based on the acquired perception data, and the vehicle’s speed and direction are continuously adjusted in real time—often via a PID controller or other conventional control algorithms to execute the parking maneuver. Although many of these algorithms have been successfully deployed in real-world vehicles, forming a mature industrial ecosystem, they remain highly dependent on accurate environmental information and precise vehicle dynamics models. As a result, any errors introduced during the perception phase can be amplified as they propagate through the planning and control stages, potentially leading to sensor deviations and the failure of the overall system.

Inspired by the Transformer model [1], this paper proposes a pure vision-based end-to-end autonomous driving approach, effectively mitigating the aforementioned problems. We sample trajectory points at certain intervals from expert trajectories and regard them as discrete tokens with positional encodings, similar to word tokens in natural language processing tasks. Subsequently, we utilize encoder and Dual-decoder structure

for autoregressive trajectory prediction, decouple the X and Y coordinates of the historical trajectories, and input them separately into the proposed model. The specific contributions of this work are as follows:

- We examined the model architectures proposed in previous end-to-end automatic parking studies and found that sequentially inputting the coordinates of x and y into a single decoder does not fully adhere to the real-world temporal order of coordinate inputs, which leads to greater trajectory errors.
- As shown in Figure 1, we propose an end-to-end neural network model **based on a dual-decoder structure**. comparative experiments demonstrate that our approach achieves more accurate and robust future trajectory predictions compared to earlier end-to-end automatic parking solutions.

II. RELATED WORK

A. BEV view in autonomous driving

Bird’s Eye View (BEV) representations offer a clear, top-down perspective, showing vehicles, roads, and obstacles at their true planar positions. Consequently, BEV is widely adopted in autonomous driving applications. However, cameras cannot directly acquire distance information for each pixel. To address this, LSS [2] compensates for the lack of depth information by predicting the depth distribution of each pixel. BEVDepth [3], on the foundation of LSS, incorporates explicit depth supervision and a Depth Refinement Module, enhancing the reliability of depth estimation. BEVDet [4] was improved on the basis of LSS and achieved higher accuracy. BEVDet4D [5], on the basis of Bevdet, introduces a time dimension, achieving the fusion of spatiotemporal features. BEVFormer [6] is a Transformer-based BEV feature generation algorithm. It can simultaneously integrate the spatial and temporal information of multi-view images and reduce the computational load, generating high-quality features. DualBEV [7] introduces a dual feature fusion strategy to effectively predict BEV probabilities. BEVfusion [8] implements a multimodal fusion perception scheme that combines LiDAR and image data, effectively preserving both geometric and semantic information. FastBEV [9] proposed Fast-Ray, thereby significantly enhancing the inference speed.

All authors are with the College of Design and Engineering, National University of Singapore, Singapore. Emails: e1373656@u.nus.edu, mpeccm@nus.edu.sg. * is the corresponding author.

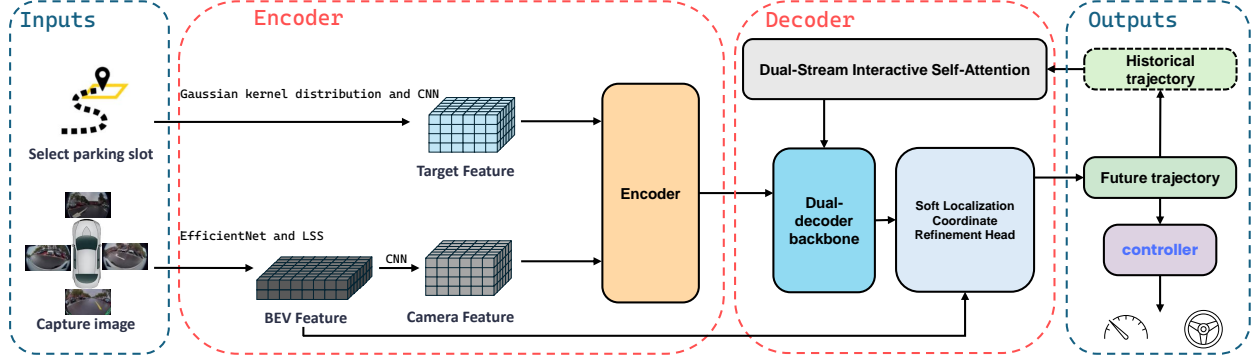


Fig. 1. Method Overview: After inputting the images captured by cameras and the coordinates of the target slot marked by the client, a Gaussian kernel distribution map is generated. After convolution, the BEV fusion feature is obtained through the attention mechanism. During the decoding process, Dual-Stream Interactive Self-Attention enables the already separated coordinates to mutually pay attention to the coordinate information of each other in the historical trajectory. Before outputting the final result, Soft Localization Coordinate Refinement Head fine-tunes the output result by utilizing the generated attention map.

B. End-to-end Autonomous Driving

In recent years, the application of neural networks in autonomous driving has significantly advanced [10]. Overall, auto-driving tasks can be broadly divided into urban driving scenarios and parking scenarios. Below, we review several representative works in each category.

Urban driving scenarios: CIL [11] directly maps front-view images, current control inputs, and high-level commands to control signals, but its performance can be sub-optimal during sharp turns. To address this issue, TCP [12] introduces a Trajectory-guided Attention mechanism, ultimately yielding multi-step control predictions based on trajectory forecasting. TransFuser [13] realizes the multimodal fusion of image and LiDAR data through the utilization of the self-attention mechanism. Transfuser++ [14] investigates the phenomenon of target point shortcuts and refines the widely used baseline model Transfuser. TAT [15] builds upon Transfuser by incorporating target points and demonstrating the significance of the target point attention mechanism. Pix2Planning [16] proposes a Transformer-based approach, in which the target position and historical trajectories are concatenated as the input sequence for the decoder. Yuandong Lyu et al. [17] fuse the bird's-eye view features based on cameras and those based on LiDAR, and ultimately output the predicted trajectory points and control signals.

Parking scenarios: Chai Runqi et al. [18] designed and verified a two-stage hybrid optimization framework that effectively combines global intelligent search and local gradient optimization. Peizhi Zhang et al. [19] proposed an end-to-end automatic parking system based on reinforcement learning, which directly maps parking space information to the optimal control instructions of the steering wheel angle. Nevertheless, it is hard to guarantee that the reward function is optimal. ParkPredict [20] introduces a nested CNN-LSTM model for

predicting trajectory points. Building on this framework, ParkPredict+ [21] replaces the LSTM with a transformer and proposes a CNN-based intention prediction network to forecast future trajectory points. Chai Runqi et al. [22] designed a motion planner based on RDNN and a tracking controller based on ALNN. E2E Parking [23] employs a transformer-based architecture that directly maps images and motion states to parking control signals once a user selects a parking spot. ParkingE2E [24] extends the E2E Parking framework. In this approach, a binary map of the target slot is convolved to generate a query vector, which is subsequently fused with the BEV feature map to predict the target point's coordinates. This method has been validated through real-world testing and has demonstrated excellent performance.

III. METHODOLOGY

A. Problem Description and Optimization Goals

We conduct training via an end-to-end neural network and define the dataset as:

$$\mathcal{D} = \{(I_{i,j}^k, X_{i,j}, Y_{i,j}, S_i, G_{i,j})\}, \quad (1)$$

where trajectory index $i \in [1, M]$, trajectory points index $j \in [1, N_i]$, camera index $k \in [1, R]$, RGB image I , ground truth trajectory point X and Y , 2D target slot coordinates S , The BEV weight probability map G based on Gaussian kernel distribution generated by the ground truth trajectory point X and Y . Reorganize the dataset as:

$$\mathcal{X}'_{i,j} = \{X_{i,\min(j+b, N_i)}\}_{b=1,2,\dots,Q}, \quad (2)$$

$$\mathcal{Y}'_{i,j} = \{Y_{i,\min(j+b, N_i)}\}_{b=1,2,\dots,Q}, \quad (3)$$

$$\mathcal{M}_{i,j} = \{G_{i,\min(j+b, N_i)}\}_{b=1,2,\dots,Q}, \quad (4)$$

and

$$\mathcal{D}' = \{(I_{i,j}^k, \mathcal{X}'_{i,j}, \mathcal{Y}'_{i,j}, S_i, \mathcal{M}_{i,j})\}, \quad (5)$$

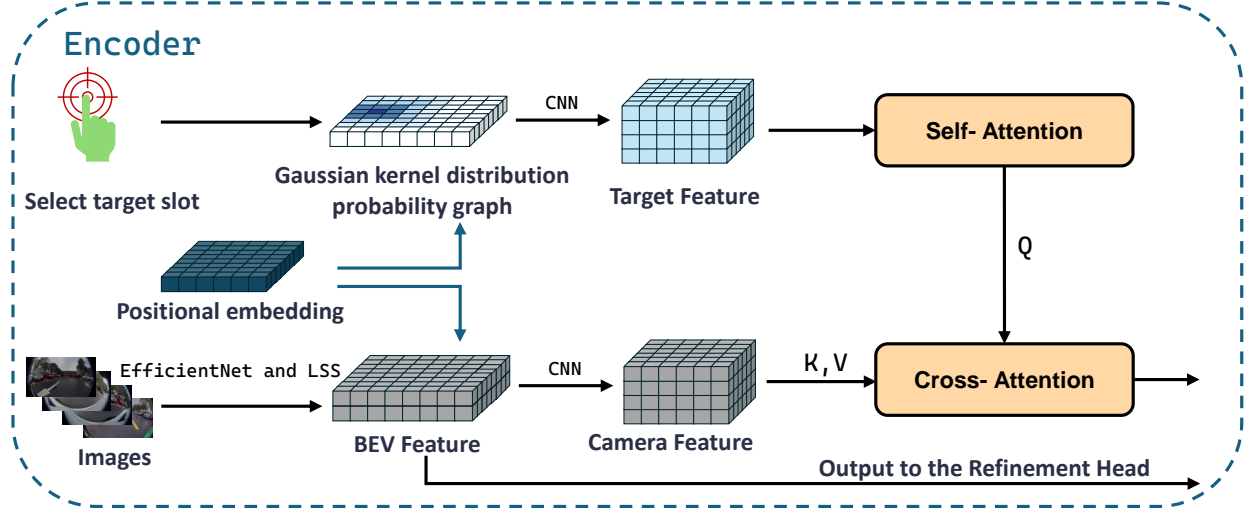


Fig. 2. The specific implementation of the Encoder structure

where Q denotes the length of the predicted trajectory points.

The optimization goals for the end-to-end network are as follows:

$$\theta^* = \arg \min_{\theta} \mathbb{E}_{(I, \mathcal{X}', \mathcal{Y}', S, \mathcal{M}) \sim \mathcal{D}'} [\mathcal{L}_x(\mathcal{X}', \mathcal{N}_{\theta}^x(I, S)) + \mathcal{L}_y(\mathcal{Y}', \mathcal{N}_{\theta}^y(I, S)) + \mathcal{L}_m(\mathcal{M}, \mathcal{N}_{\theta}^p(I, S))], \quad (6)$$

where \mathcal{L} denotes the loss function, $\mathcal{N}_{\theta}^x(I, S)$, $\mathcal{N}_{\theta}^y(I, S)$ and $\mathcal{N}_{\theta}^p(I, S)$ respectively denote the predicted X-coordinates, Y-coordinates and two-dimensional weight probability map of the output of the model.

B. End-to-End Neural Network Model Based on Camera

1) Encoder:

Query: In the ParkingE2E [24], the authors first generate a binary map with the same spatial dimensions as the BEV feature map and a depth of 1. In this map, the target slot and its surrounding fixed region are marked with a value of 1, while the area outside this range is marked with 0. After applying convolutional neural network, this map is used as a query vector to compute attention with the BEV feature map, serving as the final output of the encoder. This approach offers us significant inspiration and insight. The method of ParkingE2E clearly indicates to the model which areas need attention. However, a drawback is that it does not distinguish subregions within the target area, and the abrupt boundary transition might adversely affect continuity in attention. Therefore, in our work, we aim to give greater weight to the center while retaining some information from the surrounding areas, and this weight gradually weakens as the distance increases. We can select the desired parking slot by using the interactive tool ‘2D Nav Goal’ of Rviz and send the coordinates

to the model. Subsequently, a zero tensor with the same width and height as the BEV feature map and a depth of 1 is created. In this channel, a probability map following a Gaussian kernel distribution is generated within the specified range centered on the coordinates of target slot. Subsequently, the tensor undergoes feature extraction via same convolutional neural network, transforming it into a high-dimensional vector. The formula is as follows:

$$G(i, j) = \exp \left(-\frac{(i - cx)^2 + (j - cy)^2}{2\sigma^2} \right), \quad (7)$$

where cx, cy represent the pixel coordinates of the target slot, i, j represents the pixel coordinates on the BEV image where the Gaussian values need to be calculated. σ represents the standard deviation of the Gaussian distribution.

As shown in Table III, our ablation experiments demonstrate that the proposed method not only ensures strict selection of the target area but also enhances the granularity within the target area. This clear signal is beneficial for subsequent network training and back-propagation, thereby promoting model convergence and feature strengthening.

Key and Value: We employ EfficientNet [25] to extract image features $F_{img} \in \mathbb{R}^{C \times H_{img} \times W_{img}}$ from the RGB input, and then use the LSS [2] method to generate the BEV feature map. First, we predict the depth distribution for each pixel and map it to the 3D space. Then, we fuse the depth information with the image features. Finally, we project the fused features onto the BEV voxel grid using the camera extrinsics and intrinsics to generate the camera features. For the obtained feature map, the range in the x -direction is denoted as $[-R_x, R_x]$ m, where m indicates meters, and

the range in the y -direction is denoted as $[-R_y, R_y]$ m.

Final Output: We first conduct self-attention computation on the high-dimensional vector Query. Then, we use the BEV feature map as Key and Value. Finally, cross-attention calculation is performed among Query, Key, and Value to obtain the fused feature vector. The overall calculation process is illustrated in Figure 2.

2) Decoder: In this paper, we treat trajectory prediction as a label classification task and optimize it through cross-entropy loss.

The sequence of the trajectory: In previous works, the common practice was to feed the X-coordinate and Y-coordinate sequentially into a decoder following the temporal order. Although this method is straightforward, yet a concealed issue exists. In the real world, the movement of a vehicle is associated with the simultaneous alteration of its X-coordinate and y-coordinate. For instance, Regarding Y-coordinate Y_t at the time step t , its preceding time step ought to be Y_{t-1} . If we desire to pay additional attention to the information of X-coordinate, then X_{t-1} can also be regarded as the previous time step. In any case, it should not be X_t . Therefore, treating the X-coordinate and Y-coordinate at the same moment as two separate steps is not only a logic that makes it difficult for the model to understand, but also the longer sequence generated will lead to more error accumulation.

We separate the X and Y coordinates of the model in the ParkingE2E [24], and then train and predict using two decoders with the same number of layers. At each time step, the model simultaneously outputs the X and Y coordinates, and they can only see the historical information of the same type of coordinates. Based on the findings of the comparative experiments depicted in Table I, the experimental outcome of jointly determining the most appropriate X-coordinate and Y-coordinate at one time based on the perceived surrounding environmental information surpasses the approach of splitting the X-coordinate and Y-coordinate into two time steps.

In this paper, Trajectory serialization represents trajectory points as discrete tokens. By serializing the trajectory points, the position regression can be turned into token prediction. Subsequently, we can leverage a transformer decoder to predict the trajectory point $(p_{i,j}^x, p_{i,j}^y)$ in the ego vehicle's coordinate system. We utilize the following serialization method:

$$\text{Ser}(p_{i,j}^x) = \left\lfloor \frac{p_{i,j}^x + R_x}{2R_x} \right\rfloor \times N_t, \quad (8)$$

$$\text{Ser}(p_{i,j}^y) = \left\lfloor \frac{p_{i,j}^y + R_y}{2R_y} \right\rfloor \times N_t, \quad (9)$$

where N_t represents the maximum value that can be encoded by a token in the sequence and the symbol for

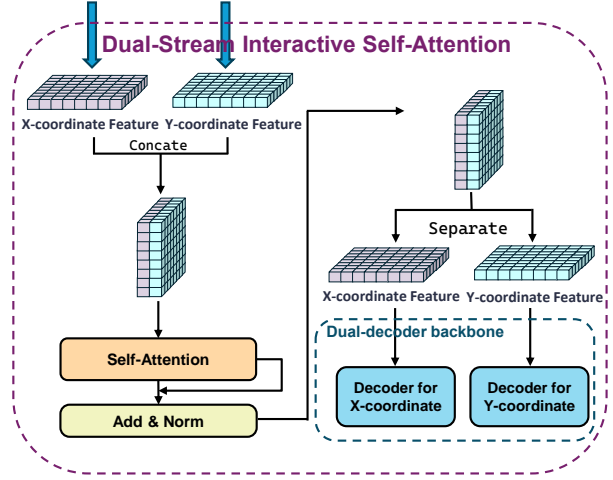


Fig. 3. The implementation logic of the Dual-Stream Interactive Self-Attention structure. The feature vectors of the X and Y coordinates are concatenated at each time step, that is, there are two vectors at each time step. After the attention calculation is completed, it is reshaped to the original size. The Batch dimension is ignored in the picture.

serializing trajectory points is denoted as $\text{Ser}(\cdot)$. R_x and R_y represent the maximum values of the predicted range in the x and y directions, respectively.

Ultimately, the trajectory sequences of the X and Y coordinates fed into the decoder are respectively:

$$[\text{BOS}, \text{Ser}(P_{i,1}^x), \text{Ser}(P_{i,2}^x), \dots, \text{Ser}(P_{i,N_i}^x), \text{EOS}], \quad (10)$$

$$[\text{BOS}, \text{Ser}(P_{i,1}^y), \text{Ser}(P_{i,2}^y), \dots, \text{Ser}(P_{i,N_i}^y), \text{EOS}], \quad (11)$$

where BOS represents the start flag and EOS represents the end flag.

TABLE I
EVALUATION OF COMPARATIVE PERFORMANCE

| Method | Haus. Dis. (m) ↓ | L2 Dis. (m) ↓ | Four. Diff. |
|------------------|------------------|---------------|--------------|
| ParkingE2E-Split | 0.2967 | 0.0808 | 1.904 |
| ParkingE2E [24] | 0.386 | 0.1193 | 2.431 |

Decoder backbone: The BEV features serve as the key and value, while the serialization sequence is used as the query. The hidden states of the X and Y coordinates of the future trajectory are generated in an autoregressive method by the dual-decoder architecture. During training, we add position embedding to the sequence points and use a mask matrix to achieve parallelization.

Dual-Stream Interactive Self-Attention: For trajectory prediction models, the X and Y coordinates usually have a close correlation. However, if only a dual-decoder structure is used for token prediction of trajectory points, it is still difficult to capture the coupling relationship

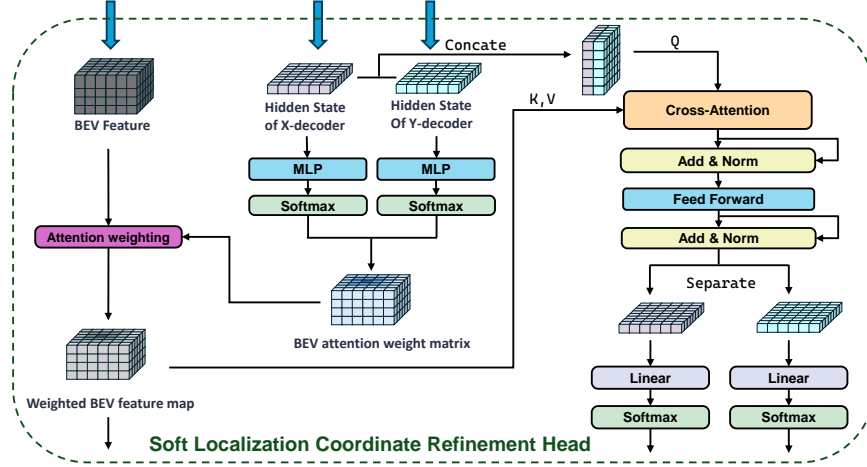


Fig. 4. The structure of the Soft Localization Coordinate Refinement Head. The BEV feature comes from the result of LSS. The hidden state vectors of the decoders of X and Y jointly predict the attention weight map of the possible area of the next trajectory point, which is multiplied by the BEV feature to generate the BEV feature weight map. The hidden state vectors interact with it through the cross-attention mechanism to achieve fine-tuning.

between the X-coordinate and Y-coordinate. In such circumstances, during the self-attention calculation, X_t coordinate only attends to the X-coordinate from X_0 to X_{t-1} , ignoring Y-coordinate features in the same period. Likewise, Y_t cannot utilize from X_0 to X_{t-1} .

To address the deficiency that the X and Y coordinates have no knowledge of each other's historical information, we propose to initially concatenate the feature vectors of the X and Y coordinates at each time step, then conduct a self-attention mechanism calculation and obtain the fused feature vector. Subsequently, the fused feature vector is separated into x and y, which are respectively input into their respective decoders. The specific structure is presented as shown in the figure 3. This attention module ensures that the high-dimensional feature vectors of the X and Y coordinate at the same time step are capable of integrating each other's feature information. It guarantees that, in the subsequent mask self-attention layers in their respective decoders, The X_t and Y_t can indirectly take into account the X-coordinate and Y-coordinate information at the 0 to $t-1$ moment.

Soft Localization Coordinate Refinement Head:

A considerable number of previous related endeavors have lacked attention to the environmental information surrounding the future trajectory points predicted by vehicles. Hence, we expect that before the model computes the final position of the next trajectory point, it should initially inspect the environmental characteristics of this area and make appropriate fine-tuning based on the environmental feature vector, and subsequently output the ultimate coordinates. Benefiting from the dual-decoder structure, We can predict the X and Y coordinates at the same time step simultaneously, which furnishes the

requisite conditions to achieve this objective.

Specifically, for the hidden state vectors of the X and Y coordinates output by the dual-decoder, after passing through the MLP layer, a learnable BEV attention weight matrix of the same size as the BEV feature map is predicted. This matrix has a predicted center point, representing the approximate position of the next trajectory point. The closer to this central point, the higher the weight. As the distance increases, the weight decreases, typically following a Gaussian kernel distribution. The specific formula is:

$$G(i, j) = \frac{\exp\left(-\frac{(i-px)^2 + (j-py)^2}{2\sigma^2}\right)}{\sum_{u=0}^{h-1} \sum_{v=0}^{w-1} \exp\left(-\frac{(u-px)^2 + (v-py)^2}{2\sigma^2}\right)}, \quad (12)$$

where px, py represent the coordinates of the Central pixel point, i, j represents the pixel coordinates on the BEV image where the Gaussian values need to be calculated. σ represents the standard deviation of the Gaussian distribution, h, w represent the height and width of the BEV image.

Then, the element-wise multiplication of the BEV attention weight matrix and the BEV feature map yields the Weighted BEV feature map, which serves as the key and value for the cross-attention calculation with the hidden state vectors of the X and Y coordinates at the current time step. Subsequently, after undergoing Feedforward layer processing, residual connection, and normalization, the final coordinates are output. In this manner, at each time step, the model is prompted to contemplate how to output more precise coordinates based on the environmental features of the region, thereby preventing the model from the issue of farsightedness.



Fig. 5. The inference results in indoor and outdoor scenarios. The orange line represents the inference trajectory of the method proposed in this paper, and the purple line represents that of ParkingE2E [24].

Note that the task of the Weighted BEV feature map is to provide environmental features near the trajectory points for the X and Y coordinates in the cross-attention calculation. It will also be output eventually, with the aim of participating in the calculation of the loss function during the training phase. The specific steps are shown in the figure 4.

C. Lateral control and longitudinal control

In our work, we adopt the similar approach as ParkingE2E [24] for Lateral and Longitudinal Control of non-omnidirectional vehicle. During the control process, we mark the initial moment of parking as t_0 . The predicted path $\mathcal{X}_{t_0}, \mathcal{Y}_{t_0} = \mathcal{N}_\theta(I_{t_0}, S)$ obtained at the initial moment is used. The relative pose from the initial moment t_0 to the current moment t can be obtained through the positioning system, denoted as $ego_{t_0 \rightarrow t}$. The target steering angle \mathcal{A}^{tar} is obtained using the RWF (Rear-wheel Feedback) method, and its expression is as follows:

$$\mathcal{A}_t^{tar} = \text{RWF}(\mathcal{X}_{t_0}, \mathcal{Y}_{t_0}, ego_{t_0 \rightarrow t}). \quad (13)$$

We achieve lateral and longitudinal control by employing a cascaded PID controller in accordance with the speed feedback \mathcal{V}^{feed} and steering feedback \mathcal{A}^{feed} from the chassis, as well as the target speed \mathcal{V}^{tar} from the setting and the target steering \mathcal{A}^{tar} from the calculation.

In this approach, the controller determines how the vehicle moves based on the future trajectory, vehicle posture and feedback signals, and only needs to focus on how to plan the route from this new starting point to the target position, without the need to constantly monitor the global positioning of the vehicle.

IV. EXPERIMENTS

A. The construction of the dataset

The dataset of this work originates from the publicly available real-vehicle dataset offered by ParkingE2E [24], which does not contain sudden obstacles and

negative samples. This dataset acquires RGB images around the vehicle through on-board cameras and employs sensor data fusion algorithms to achieve precise vehicle positioning. Meanwhile, it provides the two-dimensional coordinates of the parking target slot in the ego coordinate system for each time step.

B. Evaluation Metrics

We evaluate the inference capacity of the model by comparing the disparity between the predicted trajectory and the expert trajectory. Specifically, the following evaluation metrics are adopted:

L2 Distance: gauges the straight-line distance between the predicted trajectory and the ground-truth trajectory in the Euclidean space.

Hausdorff Distance: For each point in the predicted trajectory and the ground-truth trajectory, calculate the distance to the nearest point in the other trajectory, and take the maximum value of the two sets of minimum distances.

Fourier Descriptor Difference: transforms each of the two trajectories into the frequency domain, serving to measure the disparity between the predicted trajectory and the ground-truth trajectory.

C. comparative experiment

To guarantee the fairness of the comparative experiments, the basic training parameters employed by all the models remain consistent. The size of the BEV Features is 200×200 , The actual spatial range is $x \in [-10\text{m}, 10\text{m}]$, $y \in [-10\text{m}, 10\text{m}]$ with a resolution of 0.1 meters. In decoders, the maximum value of the trajectory serialization N_t is 1200. The trajectory decoder generates a sequence of predictions with a length of 30. Our training equipment: One NVIDIA GeForce RTX 4090 GPU. The experimental outcomes of the ParkingE2E model [24] proposed in previous works, the transfuser model (based on the GRU decoder) [13], and our proposed model are presented in Table II. The data indicates that, under the same dataset and parameters, the model we proposed has the highest accuracy rate, and all errors are reduced by **approximately 50%** compared to ParkingE2E.

TABLE II
EVALUATION OF COMPARATIVE PERFORMANCE

| Method | Haus. Dis. (m) ↓ | L2 Dis. (m) ↓ | Four. Diff. |
|---------------------|------------------|----------------|--------------|
| Our approach | 0.2156 | 0.06296 | 1.239 |
| ParkingE2E [24] | 0.386 | 0.1193 | 2.431 |
| Transfuser [13] | 0.5293 | 0.3517 | 8.674 |

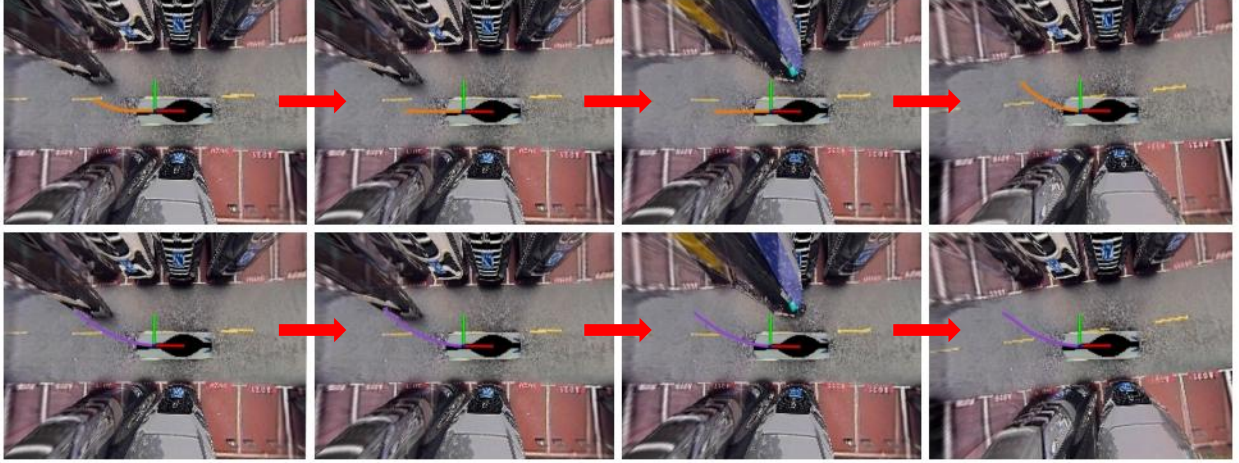


Fig. 6. The performance of the two algorithms in the same scenario where an obstacle suddenly appears. The orange line represents the inference trajectory of the method proposed in this paper, and the purple line represents that of ParkingE2E [24]. Each row of images shows the trajectory inference results of four time steps from left to right.

D. Visualization and Analysis of inference

We visualize the trajectory of model inference in Rviz, as shown in Figure 5. We also demonstrate that in the absence of special scenarios such as the sudden appearance of fast-moving objects in the training set, if other vehicles or pedestrians suddenly approach rapidly during the parking process, the model will immediately change its trajectory to prevent a collision. However, the model of ParkingE2E [24] mainly focuses its attention on the target area, thus being unable to modify the vehicle’s driving trajectory in a timely manner, as shown in Figure 6. Meanwhile, we noticed that during the parking process, if the vehicle is still at a distance from the parking space and there are no obstacles, our model will first plan a relatively straight trajectory to avoid turning too early and crossing the parking line. Only when the vehicle approaches the parking space will it start to plan a curve to achieve the turn. Moreover, with the same number of Predicted token, the trajectory inferred by ParkingE2E tends to have a larger range and sometimes deviates from the route, inferring a tortuous path. The specific effect is shown in Figure 7.

We analyzed the reason and found that the Soft Localization Coordinate Refinement Head in our model further constrains the local range through the attention heat map, making the decoded sequence more precise and convergent. Without this module, the model has less guidance or inhibition, higher degrees of freedom, and is more likely to lead to the sequence generation not converging to a shorter or more accurate trajectory. In autoregressive decoding, if the output range of the previous moment is too large or not precise enough, the subsequent predictions will accumulate errors, thereby causing this effect.

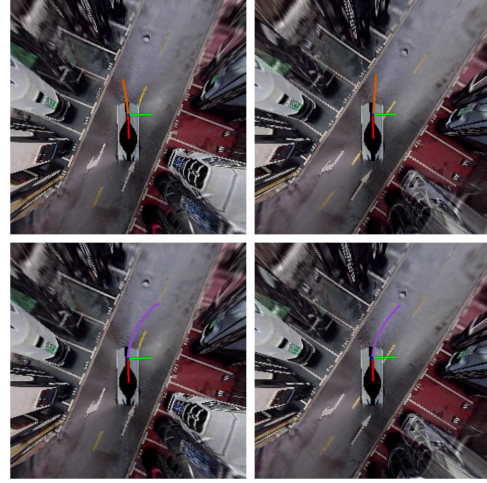


Fig. 7. The performance of the two algorithms in the same indoor scene. The orange line represents the inference trajectory of the method proposed in this paper, and the purple line represents that of ParkingE2E [24].

E. Ablation experiments

We conducted ablation experiments to verify the impact of each module. As shown in Table III, we compared our complete model, the Decoder-Binary: dual-decoder structure using binary map target queries, the Decoder-Gaussian: dual-decoder structure using Gaussian target queries, the Decoder-Gaussian-Dual: dual-decoder structure using Gaussian target queries and adding Dual-Stream Interactive Self-Attention and the Decoder-Gaussian-Refinement: dual-decoder structure using Gaussian target queries and adding Soft Localization Coordinate Refinement Head. The experiment demonstrate that each module plays a certain role and

our complete model achieves the best performance.

TABLE III
ABLATION EXPERIMENTS WITHIN THE MODEL

| Method | Haus. Dis. (m) ↓ | L2 Dis. (m) ↓ | Four. Diff. |
|-----------------------------|------------------|----------------|--------------|
| complete model | 0.2156 | 0.06296 | 1.239 |
| Decoder-Binary | 0.2967 | 0.08084 | 1.904 |
| Decoder-Gaussian | 0.2887 | 0.077 | 1.717 |
| Decoder-Gaussian-Dual | 0.2832 | 0.07556 | 1.665 |
| Decoder-Gaussian-Refinement | 0.2835 | 0.07092 | 1.631 |

V. CONCLUSIONS

In this paper, we identified the temporal issue in existing transformer-based trajectory prediction methods and addressed it by developing a dual-decoder architecture. Our experiments demonstrate that this approach improves accuracy and robustness. However, due to the limited scale and scenarios of the dataset, the current training results still have a gap compared with traditional methods. Therefore, we have only conducted trajectory inference on images collected from real vehicles so far and have not yet carried out road tests. In future work, we plan to expand the dataset, strengthen the model's capability of handling more complex tasks, and carry out real vehicle tests. The current query of encoder still determines a fixed range by parameters. In the future, we can make an attempt to employ reinforcement learning to allow the model to focus on a target range more flexibly.

Overall, we are convinced that the learning-based trajectory prediction is capable of continuous evolution, genuinely attaining and exceeding the level of traditional methods, and bringing a wider range of development space to the domain of autonomous driving.

REFERENCES

- [1] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [2] J. Philion and S. Fidler, "Lift, splat, shoot: Encoding images from arbitrary camera rigs by implicitly unprojecting to 3d," in *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIV 16*, pp. 194–210, Springer, 2020.
- [3] Y. Li, Z. Ge, G. Yu, J. Yang, Z. Wang, Y. Shi, J. Sun, and Z. Li, "Bevdepth: Acquisition of reliable depth for multi-view 3d object detection," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 37, pp. 1477–1485, 2023.
- [4] J. Huang, G. Huang, Z. Zhu, Y. Ye, and D. Du, "Bevdet: High-performance multi-camera 3d object detection in bird-eye-view," *arXiv preprint arXiv:2112.11790*, 2021.
- [5] J. Huang and G. Huang, "Bevdet4d: Exploit temporal cues in multi-camera 3d object detection," *arXiv preprint arXiv:2203.17054*, 2022.
- [6] Z. Li, W. Wang, H. Li, E. Xie, C. Sima, T. Lu, Q. Yu, and J. Dai, "Bevformer: Learning bird's-eye-view representation from multi-camera images via spatiotemporal transformers," *arXiv e-prints*, pp. arXiv:2203.17054, 2022.
- [7] P. Li, W. Shen, Q. Huang, and D. Cui, "Dualbev: Unifying dual view transformation with probabilistic correspondences," in *European Conference on Computer Vision*, pp. 286–302, Springer, 2024.
- [8] Z. Liu, H. Tang, A. Amini, X. Yang, H. Mao, D. L. Rus, and S. Han, "Bevfusion: Multi-task multi-sensor fusion with unified bird's-eye view representation," in *2023 IEEE international conference on robotics and automation (ICRA)*, pp. 2774–2781, IEEE, 2023.
- [9] Y. Li, B. Huang, Z. Chen, Y. Cui, F. Liang, M. Shen, F. Liu, E. Xie, L. Sheng, W. Ouyang, *et al.*, "Fast-bev: A fast and strong bird's-eye view perception baseline," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- [10] M. Bojarski, D. Del Testa, D. Dworakowski, B. Firner, B. Flepp, P. Goyal, L. D. Jackel, M. Monfort, U. Muller, J. Zhang, *et al.*, "End to end learning for self-driving cars," *arXiv preprint arXiv:1604.07316*, 2016.
- [11] F. Codevilla, M. Müller, A. López, V. Koltun, and A. Dosovitskiy, "End-to-end driving via conditional imitation learning," in *2018 IEEE international conference on robotics and automation (ICRA)*, pp. 4693–4700, IEEE, 2018.
- [12] P. Wu, X. Jia, L. Chen, J. Yan, H. Li, and Y. Qiao, "Trajectory-guided control prediction for end-to-end autonomous driving: A simple yet strong baseline," *Advances in Neural Information Processing Systems*, vol. 35, pp. 6119–6132, 2022.
- [13] K. Chitta, A. Prakash, B. Jaeger, Z. Yu, K. Renz, and A. Geiger, "Transfuser: Imitation with transformer-based sensor fusion for autonomous driving," *IEEE transactions on pattern analysis and machine intelligence*, vol. 45, no. 11, pp. 12878–12895, 2022.
- [14] B. Jaeger, K. Chitta, and A. Geiger, "Hidden biases of end-to-end driving models," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 8240–8249, 2023.
- [15] A. Shamsoshoara, S. B. Salih, and P. Aghazadeh, "Swaptransformer: highway overtaking tactical planner model via imitation learning on osha dataset," *IEEE Access*, 2024.
- [16] X. Mu, T. Qin, S. Zhang, C. Xu, and M. Yang, "Pix2planning: End-to-end planning by vision-language model for autonomous driving on carla simulator," in *2024 IEEE Intelligent Vehicles Symposium (IV)*, pp. 2383–2390, IEEE, 2024.
- [17] Y. Lyu, X. Tan, Z. Yu, and Z. Fan, "Sensor fusion and motion planning with unified bird's-eye view representation for end-to-end autonomous driving," in *2024 International Joint Conference on Neural Networks (IJCNN)*, pp. 1–8, IEEE, 2024.
- [18] R. Chai, A. Tsourdos, A. Savvaris, S. Chai, and Y. Xia, "Two-stage trajectory optimization for autonomous ground vehicles parking maneuver," *IEEE Transactions on Industrial Informatics*, vol. 15, no. 7, pp. 3899–3909, 2018.
- [19] P. Zhang, L. Xiong, Z. Yu, P. Fang, S. Yan, J. Yao, and Y. Zhou, "Reinforcement learning-based end-to-end parking for automatic parking system," *Sensors*, vol. 19, no. 18, p. 3996, 2019.
- [20] X. Shen, I. Batkovic, V. Govindarajan, P. Falcone, T. Darrell, and F. Borrelli, "Parkpredict: Motion and intent prediction of vehicles in parking lots," in *2020 IEEE Intelligent Vehicles Symposium (IV)*, pp. 1170–1175, IEEE, 2020.
- [21] X. Shen, M. Lacayo, N. Guggilla, and F. Borrelli, "Parkpredict+: Multimodal intent and motion prediction for vehicles in parking lots with cnn and transformer," in *2022 IEEE 25th International Conference on Intelligent Transportation Systems (ITSC)*, pp. 3999–4004, IEEE, 2022.
- [22] R. Chai, D. Liu, T. Liu, A. Tsourdos, Y. Xia, and S. Chai, "Deep learning-based trajectory planning and control for autonomous ground vehicle parking maneuver," *IEEE Transactions on Automation Science and Engineering*, vol. 20, no. 3, pp. 1633–1647, 2022.
- [23] Y. Yang, D. Chen, T. Qin, X. Mu, C. Xu, and M. Yang, "E2e parking: Autonomous parking by the end-to-end neural network on the carla simulator," in *2024 IEEE Intelligent Vehicles Symposium (IV)*, pp. 2375–2382, IEEE, 2024.
- [24] C. Li, Z. Ji, Z. Chen, T. Qin, and M. Yang, "Parking2e: Camera-based end-to-end parking network, from images to planning," in *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 13206–13212, IEEE, 2024.
- [25] M. Tan and Q. Le, "Efficientnet: Rethinking model scaling for convolutional neural networks," in *International conference on machine learning*, pp. 6105–6114, PMLR, 2019.