

Divide and Merge: Motion and Semantic Learning in End-to-End Autonomous Driving

Anonymous authors

Paper under double-blind review

Abstract

Perceiving the environment and its changes over time corresponds to two fundamental yet heterogeneous types of information: semantics and motion. Previous end-to-end autonomous driving works represent both types of information in a single feature vector. However, including motion related tasks, such as prediction and planning, impairs detection and tracking performance, a phenomenon known as negative transfer in multi-task learning. To address this issue, we propose Neural-Bayes motion decoding, a novel parallel detection, tracking, and prediction method that separates semantic and motion learning. Specifically, we employ a set of learned motion queries that operate in parallel with detection and tracking queries, sharing a unified set of recursively updated reference points. Moreover, we employ interactive semantic decoding to enhance information exchange in semantic tasks, promoting positive transfer. Experiments on the nuScenes dataset with UniAD and SparseDrive confirm the effectiveness of our divide and merge approach, resulting in performance improvements across perception, prediction, and planning. The code will be released.

1 Introduction

Modular end-to-end (E2E) autonomous driving (AD) is gaining attention for combining the strengths of traditional pipeline methods with strict E2E approaches. In this framework, perception, prediction, and planning form the core set of tasks, which ideally complement one another to enhance overall system performance. However, the modular E2E framework also presents a multi-task learning challenge. A poorly designed multi-task learning structure could not only fail to facilitate mutual learning but also adversely affect individual tasks, a phenomenon known as negative transfer (Crawshaw, 2020). The prevalent modular E2E approaches (Hu et al., 2023; Jiang et al., 2023; Zheng et al., 2025; Sun et al., 2024) typically employ a sequential structure (Fig. 1a). This structure aligns with how humans perform driving tasks and has demonstrated promising planning performance. However, these approaches exhibit negative transfer in object detection and tracking. In other words, the perception performance of jointly trained E2E models is typically inferior to those trained without the motion prediction and planning tasks.

We analyze the underlying causes of negative transfer by inspecting the types of learned heterogeneous information: semantic and motion. Semantic information encompasses the categories of surrounding objects, lanes, crossings, *etc.*, while motion information describes the temporal changes occurring within the environment. Sequential methods (Hu et al., 2023; Jiang et al., 2023; Zheng et al., 2025; Doll et al., 2024) execute these two processes in succession. They first conduct detection and tracking and then use the extracted object features for trajectory prediction. This sequential design forces the features to contain motion information, compromising the initially learned semantic and leading to negative transfer in perception. The SHAP values analysis (Lundberg & Lee, 2017) provides supporting evidence for our argument. Another E2E structure is depicted in Fig. 1b. It executes most tasks with different heads in parallel, as PARA-Drive (Weng et al., 2024) and NMP (Zeng et al., 2019). However, since detection and prediction remain sequential, the issue of negative transfer persists.

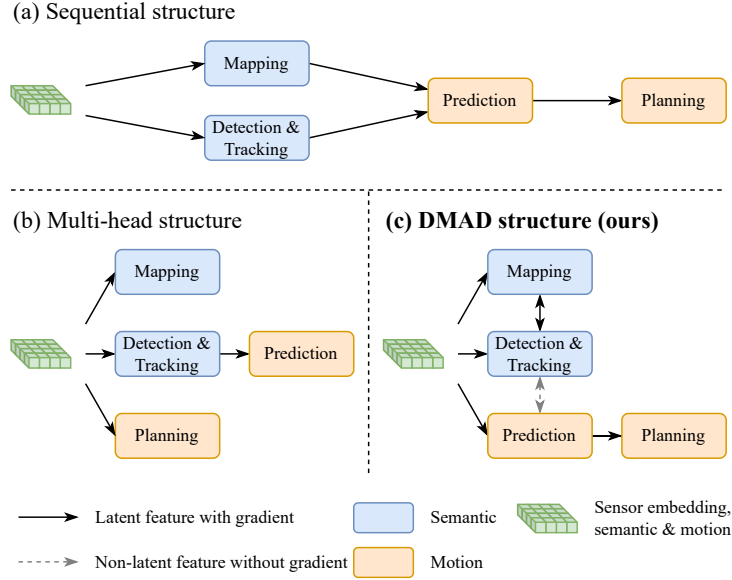


Figure 1: **Comparison of E2E structures.** In (a), semantic and motion learning occur sequentially. In (b), the multi-head structure parallelizes tasks with different heads; however, motion and semantic learning remain sequential in detection, tracking, and prediction. In (c), semantic and motion learning are performed in parallel without latent feature sharing or gradient propagation. In contrast, the exchange of information between the object and map perception modules is enhanced.

In this work, we propose **DMAD structure** (Fig. 1c), **D**ividing and **M**erging motion and semantic learning for E2E **A**utonomous **D**riving. DMAD addresses the issue of negative transfer by separating semantic and motion learning. Furthermore, it leverages correlations among semantic tasks by merging them.

For dividing, we propose **Neural-Bayes motion decoder**. We maintain a set of motion queries that attend to the sensor embeddings parallel to the object (detection and tracking) queries. The key difference between motion and object queries is that they are decoded into past and future trajectories rather than bounding boxes with classes. Motion and object queries share a single set of reference points, updated recursively by detection and prediction. It allows only limited information exchange between both types of queries, mediated through the reference points without gradient flow. Moreover, we calculate the object’s velocity using the predicted trajectory with finite differences, thereby removing the requirement for object queries to learn the velocity directly. In this manner, the object query focuses on learning semantic and appearance features, while the motion query is dedicated to capturing motion features. The two types of heterogeneous information are learned separately along distinct paths, effectively preventing negative transfer. Notably, the DMAD structure promotes motion learning to the same level of semantic learning, treating detection, tracking, and prediction as concurrent tasks for the first time, to the best of our knowledge.

For merging, we propose **interactive semantic decoder** to enhance the exchange of semantic insights in detection and map segmentation. Object perception and map perception are inherently related tasks. Previous methods often overlook this connection, typically executing the two along parallel paths (Hu et al., 2023; Jiang et al., 2023; Zheng et al., 2025). DualAD (Doll et al., 2024) leverages this correlation but allows only object perception to learn from the map. Our method uses layer-wise iterative self-attention (Vaswani, 2017) to enable mutual learning between object and map tasks, fostering positive transfer.

Experiments on the nuScenes (Caesar et al., 2020) dataset showcase the effectiveness of DMAD structure in mitigating negative transfer. Our approach achieves significant performance gains in perception and prediction, which benefits the planning module and outperforms state-of-the-art (SOTA) E2E AD models.

Our key contributions are summarized as follows:

- We examine the similarity and heterogeneity among tasks in modular E2E AD and argue that the prevailing design—learning information for conflicting tasks within a single feature—is the cause of negative transfer in perception. We analyze SHAP values to validate this hypothesis. Conversely, we propose that information exchange between similar tasks can facilitate positive transfer.
- We propose DMAD, a modular E2E AD paradigm that divides and merges tasks according to the information they are supposed to learn. This design eliminates negative transfer between different types of tasks while reinforcing positive transfer among similar tasks.
- We introduce two decoders: the Neural-Bayes motion decoder for concurrent trajectory prediction with object detection and tracking; the interactive semantic decoder to enhance information sharing between object and map perception. The proposed decoders improve existing SOTA methods, leading to better performance across all tasks.

2 Related Work

Semantic learning. Semantic learning includes object detection and map segmentation. Multi-view cameras have become popular due to their cost-effectiveness and strong capability in capturing semantic information. Current SOTA object detection and mapping approaches are built on the DETR (Carion et al., 2020) architecture, utilizing a set of queries to extract semantic information from environment features through cross-attention (Vaswani, 2017) mechanisms. Sparse methods (Wang et al., 2022; Lin et al., 2022) learn semantic information by projecting queries onto the corresponding image features, focusing on the relevant regions. The PETR series (Liu et al., 2022; 2023; Wang et al., 2023) embed 3D positional encoding directly into 2D image features, eliminating the need for query projection. Another line of work aggregates all image features into a bird’s-eye view (BEV) feature (Phillion & Fidler, 2020; Li et al., 2022; Yang et al., 2023; Pan et al., 2024; Liao et al., 2023; 2024). Propagating the object queries over time enables multi-object tracking (Zeng et al., 2022; Meinhardt et al., 2022). This same technique is also used in map perception (Chen et al., 2025). Although tracking is also a motion-related task, we classify it as a semantic task, as query-based trackers learn only velocities as the motion information, which we elaborate in Appendix A.

Motion learning. By motion, we refer to trajectory prediction and planning. Trajectory prediction studies typically use the ground truth of objects’ historical trajectories along with high-definition maps as inputs. Early approaches (Chai et al., 2019; Cui et al., 2019; Bansal et al., 2019) rasterize maps and trajectories into a BEV image, using CNNs to extract scene features. Vectorized methods (Gao et al., 2020; Zhou et al., 2022) represent elements using polygons and polylines, using GNNs or Transformers to encode the scene (Ngiam et al., 2022; Wagner et al., 2024; Shi et al., 2022; Gu et al., 2021; Zhang et al., 2024).

For planning, imitation learning is a straightforward approach to E2E planning, where a neural network is trained to plan future trajectories or control signals directly from sensor data, minimizing the distance between the planned path and the expert driving policy (Bojarski, 2016; Prakash et al., 2021; Chen & Krähenbühl, 2022). Many approaches incorporate semantic tasks as auxiliary components to support E2E planning, using the nuScenes (Caesar et al., 2020) dataset and open-loop evaluation. These methods go beyond pure motion learning and are presented in the next paragraph. AD-MLP (Zhai et al., 2023) and Ego-MLP (Li et al., 2024) utilize only the ego vehicle’s past motion states and surpass methods that rely on sensor inputs in open-loop evaluation. It aligns with our argument that semantics and motion are heterogeneous: AD-MLP and Ego-MLP can concentrate on learning from expert motion data without interference by irrelevant semantic information, thereby achieving superior open-loop planning performance.

Joint semantic and motion learning. E2E perception and prediction approaches learn semantics and motion jointly. The pioneering work FaF (Luo et al., 2018) uses a prediction head, in addition to the detection head, to decode the object features into future trajectories. Some works (Casas et al., 2018; Djuric et al., 2021; Fadadu et al., 2022) enhance it with intention-based prediction and refinement. PnPNet (Liang et al., 2020) and PTP (Weng et al., 2021) involve tracking, *i.e.*, jointly optimizing detection, association, and prediction tasks. While PTP performs tracking and prediction in parallel, it cannot predict newly

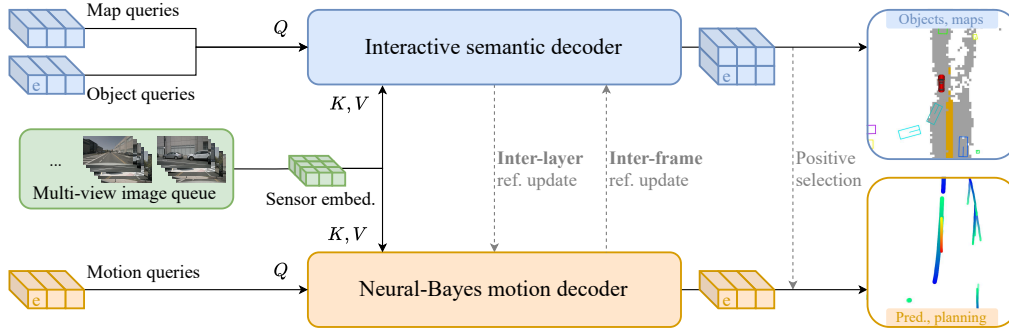


Figure 2: **An overview of DMAD.** A backbone processes multi-view images into sensor embeddings. Map and object queries are initialized, then interactively attend to the sensor embeddings for map and object perception. Motion queries, mapped one-to-one with object queries, share reference points that are iteratively updated. Finally, motion queries corresponding to detected objects are decoded into future trajectories. The ego motion query (“e”) is used for planning. Gray dashed lines indicate operations without gradient flow.

emerging objects due to the lack of concurrent detection—a limitation our method successfully overcomes. ViP3D (Gu et al., 2023) first extends the query-based detection and tracking framework (Zeng et al., 2022) to prediction. Each query represents an object and propagates across frames. In each frame, queries are decoded into bounding boxes and trajectories using high-definition maps as additional context.

To include planning, NMP (Zeng et al., 2019) extends IntentNet (Casas et al., 2018) with a sampling-based planning module, where prediction is leveraged to minimize collisions during the planning process. Other works (Chitta et al., 2021; Casas et al., 2021; Hu et al., 2022) incorporate map perception as an auxiliary task. With the growing popularity of query-based object detectors (Carion et al., 2020; Li et al., 2022) and trackers (Zeng et al., 2022; Meinhardt et al., 2022), recent modular E2E AD approaches represent objects as queries, similar to ViP3D (Gu et al., 2023). UniAD (Hu et al., 2023) and its variants (Doll et al., 2024; Weng et al., 2024) retain the query propagation mechanism for tracking, aiming to explicitly model objects’ historical motion. In contrast, VAD (Jiang et al., 2023) and GenAD (Zheng et al., 2025) do not perform tracking, predicting trajectories based on the temporal information embedded within the BEV feature. The main issue with these methods is that they attempt to use a single feature (query) to represent an object’s appearance and motion. Compared to pure semantic learning, motion occupies a portion of the feature channels but fails to contribute to perception, resulting in a negative transfer in the perception module. Our work effectively addresses this issue.

3 Method

Figure 2 shows an overview of DMAD structure. Sensor embeddings are extracted from multi-view camera images and are shared across all tasks, including detection, tracking, mapping, prediction, and planning. We initialize three distinct types of queries—object, map, and motion—which attend to the sensor embeddings to extract the specific information required for each respective task. Based on the type of information learned, the decoding process is divided into two pathways. On one way, object and map decoding are jointly performed within the **Interactive semantic decoder**, where both types of queries iteratively exchange latent semantic information at each decoding layer. On the other way, motion queries extract motion information from the sensor embeddings within the **Neural-Bayes motion decoder**. Each motion query is paired with an object query, using the object’s coordinates as a reference point at each decoding layer. After decoding each frame, the motion query’s predicted future waypoint becomes the object query’s reference point in the next frame, similar to the recursion of a Bayes filter (Thrun et al., 2005). The exchange of reference points is always without gradient. At last, the motion queries are passed on to the planning module. The system is fully E2E trainable, with motion and semantic gradients propagated in distinct paths.

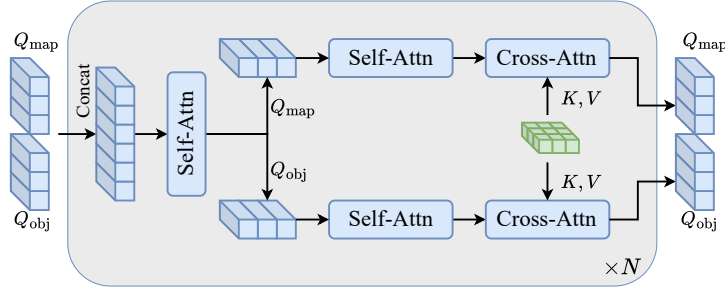


Figure 3: **Interactive semantic decoding.** Object and map queries are concatenated and interact through a self-attention module before being separated to independently attend to the sensor embeddings. This process is repeated across N stacked layers.

3.1 Interactive Semantic Decoder

To leverage the semantic correlation between individual objects and map elements, we introduce the Interactive Semantic Decoder. In contrast to the unidirectional interaction in DualAD (Doll et al., 2024), our approach enables a bidirectional exchange of information.

We initialize a set of object queries $Q_{\text{obj}} \in \mathbb{R}^{N_{\text{obj}} \times d}$ and a set of map queries $Q_{\text{map}} \in \mathbb{R}^{N_{\text{map}} \times d}$. The number of queries could be different, while the dimensions d must be the same. Each decoding layer first concatenates both types of queries. Self-attention (Vaswani, 2017) is then applied, where both tasks exchange their semantic information. Subsequently, the two types of queries are divided, each performing self-attention and cross-attention on the sensor embeddings, respectively, as shown in Fig. 3.

After interactive semantic decoding, each object query is classified into a category c and regressed into a vector $[\Delta x, \Delta y, \Delta z, w, h, l, \theta]^T$. The object query is associated with a reference point $[x_{\text{ref}}, y_{\text{ref}}, z_{\text{ref}}]^T$. Rather than directly learning the absolute coordinates of the object, it learns the offsets relative to its corresponding reference points. Thus, the bounding boxes can be represented as $[x_{\text{ref}} + \Delta x, y_{\text{ref}} + \Delta y, z_{\text{ref}} + \Delta z, w, h, l, \theta]^T$. Notably, velocities are not regressed, as they pertain to motion information. We design the object queries to focus solely on semantic information, *i.e.*, the object’s category, center point, size, and orientation.

3.2 Neural-Bayes Motion Decoder

We introduce a novel motion decoder operating in parallel with the semantic decoder, aimed at fully decoupling motion and semantic learning to reduce the negative transfer in semantic tasks. Given the correlation between motion and semantics, we design a recursive process to facilitate the exchange of human-readable information between the two decoders as illustrated in Fig. 4, which comprises the processes of prediction, measurement, and updating, similar to the Bayes filter (Thrun et al., 2005). Appendix B provides a brief introduction to the Bayes filter. We proceed with the elaboration of the proposed motion decoder.

Initialization. We initialize a set of motion queries $Q_{\text{mt}} \in \mathbb{R}^{N_{\text{mt}} \times d}$ in the same way we initialize object queries. The motion queries correspond one-to-one with the object queries, *i.e.*, $N_{\text{mt}} = N_{\text{obj}}$. However, since they do not directly interact in the latent space, their dimensionalities d can differ. Each motion query represents the motion state of an object, although the model does not initially know whether the object exists. Additionally, motion queries and object queries share a common set of reference points.

Measurement. The detection, already introduced in Sec. 3.1, is treated as the measurement in Bayes filter. After each semantic decoding layer, the object queries are regressed, yielding the coordinate vectors $\text{ref} = [x, y, z]^T$ of the tentative object, which then serves as reference points for the next layer:

$$\text{ref}^{l+1} = f_{\text{reg}}(f_{\text{Semantic-Dec}}^l(Q_{\text{obj}}^l, Z, \text{ref}^l)), \quad (1)$$

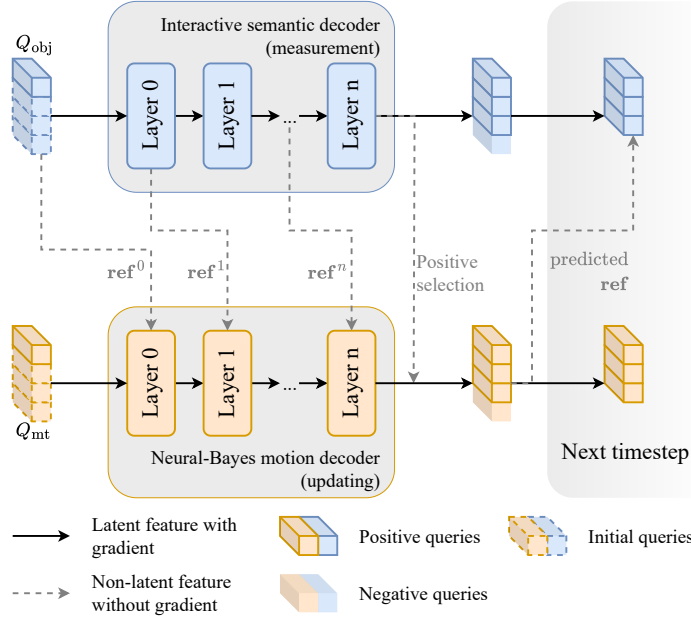


Figure 4: **Neural-Bayes motion decoding.** After each decoding layer, the semantic decoder updates the reference points, which are then shared with the motion decoder. At the end of each frame, positive object query indices are used to select corresponding motion queries and are together propagated to the subsequent frame, with the motion query predictions serving as reference points for the next frame. This process is similar to the measurement, updating, and prediction steps in a Bayes filter. Map queries, ego queries and sensor embeddings are omitted for simplicity.

where the superscript denotes the layer and Z is the sensor embeddings.

Updating. With the reference points \mathbf{ref}^l from the semantic decoding (the inter-layer reference points update in Fig. 2), the motion queries also attend to the sensor embeddings via cross-attention:

$$Q_{\text{mt}}^{l+1} = f_{\text{Motion-Dec}}^l(Q_{\text{mt}}^l, Z, \mathbf{ref}^l), \quad (2)$$

where the motion queries are updated conditioned on the measured reference points.

Prediction. We employ MLPs to extract trajectories from the motion queries. We note that motion extraction occurs in two stages: first through the unimodal trajectory construction, followed by the multimodal prediction.

The first stage computes the unimodal velocity and future reference points, guiding the motion query to learn aggregated motion states from the past and predict the near future. It produces a single trajectory that spans from the past timestep t_{past} to the future timestep $t_{\text{fut-1}}$. The velocity is calculated using the finite difference method on waypoints around the current timestep. We use the first future waypoint as the initial reference point for the object query in the next frame, *i.e.*, inter-frame reference points update in Fig. 2, for object tracking.

The second stage performs multimodal intention modeling and generates multiple future trajectories within the future $t_{\text{fut-2}}$ timesteps, along with their corresponding confidence scores.

Tracking. Multi-object tracking is performed using the query propagation mechanism (Zeng et al., 2022; Lin et al., 2023). Each object query is associated with an unique instance ID. A positive query propagates across consecutive frames, ensuring that corresponding detections are assigned the same ID. During training,

object queries associated with ground truth are referred to as positive queries; during inference, positivity is determined by whether the confidence score exceeds a specified threshold. The propagation of motion queries follows that of object queries, as they are related. This mechanism enables continuous measuring, updating, and predicting, similar to the Bayes filter.

4 Experiments

We conduct experiments on the nuScenes (Caesar et al., 2020) dataset to validate the effectiveness of our method. We present results in three parts. The first part focuses on perception (detection, tracking, and mapping). In the second part, we evaluate motion prediction and planning. Lastly, we provide an extensive ablation study and SHAP values (Lundberg & Lee, 2017) visualization.

4.1 Training Configuration

We reproduce UniAD (Hu et al., 2023) and SparseDrive (Sun et al., 2024) as baselines. Both utilize the query propagation mechanism; however, UniAD extracts dense BEV features from image inputs, while SparseDrive employs sparse scene representations. Beside the aforementioned tasks, UniAD additionally performs occupancy prediction. We also retain the occupancy module in comparisons with UniAD for task consistency. As occupancy prediction serves merely as another representation of upstream tasks, we describe it in Appendix C. We adhere as closely as possible to default configurations of the baseline; however, to ensure a rigorous comparisons, some adjustments are made. Following paragraphs outline the adjustments and the rationale behind them.

Two-stage training. We follow the two-stage training scheme of our baseline. In the first stage, we train object detection, tracking, and mapping. In the second stage, we train all modules together. Notably, because our tracking relies on reference points provided by unimodal prediction, we incorporate unimodal prediction training in the first stage. Multimodal prediction is trained only in the second stage, which is consistent with the baseline.

Queue length. Since AD is a time-dependent task, the model typically processes a sequence of consecutive frames as a training sample. The number of input frames, *i.e.*, the queue length q , defines the temporal horizon the model can capture, impacting the performance of related tasks. UniAD employs different queue lengths across its two training stages: 5 in the first stage and 3 in the second. The reduced queue length in the second stage degrades perception performance due to reduced temporal aggregation, shown in Appendix D. This degrading hinders the identification of negative transfer effects caused by the sequential structure. To mitigate this interference, we standardize the queue length to 3 across both training stages in comparisons with UniAD. Unless otherwise specified, the performance of UniAD in all result tables is reproduced with a queue length of 3 using the official codebase (UniAD-contributors, 2023). SparseDrive does not have this issue, and we use the default setting of 4.

Ego query represents the features directly used for motion planning, which is intended to capture the motion information of the ego vehicle. SparseDrive generates the ego query from the front camera image and the estimated previous ego status, which blends semantics and motion, thus contradicting our dividing design. To align with our proposal, we eliminate the use of the front image for the ego query when applying DMAD to SparseDrive. For UniAD, we retain the planning module unchanged, as it initializes the ego query randomly.

4.2 Perception

Metrics. For object detection and tracking, we use the metrics defined in the nuScenes benchmark. The primary metrics for detection are nuScenes Detection Score (NDS) and mean average precision (mAP). For multiple object tracking, we report the average multi-object tracking accuracy (AMOTA) and the average

| Method | NDS↑ | mAP↑ | mAVE↓ |
|--------------------------------|---------------|----------------------|----------------------|
| VAD (Jiang et al., 2023) | 0.460 | 0.330 | 0.405 |
| GenAD (Zheng et al., 2025) | 0.280 | 0.213 | 0.669 |
| PARA-Drive (Weng et al., 2024) | 0.480 | 0.370 | - |
| UniAD - stage 1 | 0.497 | 0.382 | 0.411 |
| UniAD - stage 2 | 0.491 (-1.2%) | 0.377 (-1.3%) | 0.412 (+0.2%) |
| DMAD - stage 1 | 0.504 | 0.395 | 0.406 |
| DMAD - stage 2 | 0.506 (+0.4%) | 0.396 (+0.3%) | 0.395 (-2.7%) |
| SparseDrive - stage 1 | 0.531 | 0.419 | 0.257 |
| SparseDrive - stage 2 | 0.523 (-1.5%) | 0.417 (-0.5%) | 0.269 (+4.7%) |
| SparseDMAD - stage 1 | 0.536 | 0.424 | 0.260 |
| SparseDMAD - stage 2 | 0.534 (-0.4%) | 0.427 (+0.7%) | 0.253 (-2.7%) |

Table 1: **Object detection results.** The performance changes in stage 2 are expressed as percentages, with **red** indicating a decline and **blue** representing improvement.

| Method | AMOTA↑ | AMOTP↓ | IDS↓ |
|--------------------------------|---------------|------------------|--------------------|
| ViP3D (Gu et al., 2023) | 0.217 | 1.63 | - |
| MUTR3D (Zhang et al., 2022) | 0.294 | 1.50 | 3822 |
| PARA-Drive (Weng et al., 2024) | 0.350 | - | - |
| UniAD - stage 1 | 0.374 | 1.31 | 816 |
| UniAD - stage 2 | 0.354 (-5.3%) | 1.34 (+2.3%) | 1381 (+69%) |
| DMAD - stage 1 | 0.394 | 1.32 | 781 |
| DMAD - stage 2 | 0.393 (-0.3%) | 1.30 (-1.5%) | 767 (-1.8%) |
| SparseDrive - stage 1 | 0.395 | 1.25 | 602 |
| SparseDrive - stage 2 | 0.376 (-4.8%) | 1.26 (+0.8%) | 559 (-7.1%) |
| SparseDMAD - stage 1 | 0.396 | 1.23 | 608 |
| SparseDMAD - stage 2 | 0.395 (-0.3%) | 1.23 (0%) | 571 (-6.1%) |

Table 2: **Multi-object tracking results.**

multi-object tracking precision (AMOTP). For map segmentation, we use the intersection over union (IoU) metric of drivable areas, lanes, and dividers. Vectorized mapping adopts mAP of lane divider, pedestrian crossing and road boundary.

Object detection. Table 1 presents the detection performance across two training stages. In the first stage, thanks to the interactive semantic decoding, our approach slightly outperforms the baseline. After the second stage of training, baseline’s performance shows a decline. In contrast, our method preserves the perceptual performance of the first stage, benefiting from separated motion learning that mitigates negative transfer. Our method finally surpasses UniAD and SparseDrive by 3.1% and 2.1% in NDS, respectively.

Multi-object tracking. Due to using a single feature vector to represent semantics and motion, UniAD and SparseDrive exhibit negative transfer of 5.3% and 4.8% in AMOTA, as shown in Tab. 2. Our dividing design enables object queries to learn about appearance more effectively. At the same time, unimodal predictions offer enhanced tracking reference points. Consequently, our method achieves a gain of 11.0% and 5.1% in AMOTA, respectively.

Map perception. UniAD does not encounter negative transfer in map segmentation. Leveraging the advantages of interactive semantic decoding, our method marginally surpasses UniAD. Our method mitigates

| Method | Lanes \uparrow | Drivable \uparrow | Dividers \uparrow |
|--------------------------------|----------------------|---------------------|----------------------|
| BEVFormer (Li et al., 2022) | 0.239 | 0.775 | - |
| PARA-Drive (Weng et al., 2024) | 0.330 | <u>0.710</u> | - |
| UniAD - stage 1 | 0.293 | 0.650 | 0.248 |
| UniAD - stage 2 | 0.312 (+6.5%) | 0.678 (+4.3%) | <u>0.267</u> (+7.7%) |
| DMAD - stage 1 | 0.292 | 0.655 | 0.242 |
| DMAD - stage 2 | <u>0.321</u> (+9.9%) | 0.691 (+5.5%) | 0.271 (+12%) |

(a) Map segmentation results.

| Method | AP _{ped} \uparrow | AP _{divider} \uparrow | AP _{boundary} \uparrow | mAP \uparrow |
|---------------------------|------------------------------|----------------------------------|-----------------------------------|----------------------|
| MapTR (Liao et al., 2023) | 0.562 | 0.598 | <u>0.601</u> | 0.587 |
| VAD (Jiang et al., 2023) | 0.406 | 0.515 | 0.506 | 0.476 |
| SparseDrive - stage 1 | 0.533 | 0.579 | 0.575 | 0.562 |
| SparseDrive - stage 2 | 0.494 (-7.3%) | 0.569 (-1.7%) | 0.583 (+1.4%) | 0.549 (-2.3%) |
| SparseDMAD - stage 1 | 0.553 | <u>0.599</u> | 0.606 | <u>0.586</u> |
| SparseDMAD - stage 2 | <u>0.554</u> (+0.2%) | 0.601 (+0.3%) | 0.606 (0%) | 0.587 (+0.2%) |

(b) Vectorized mapping results.

Table 3: Map perception results.

| Method | EPA \uparrow | | minADE \downarrow | |
|----------------------------|----------------|--------------|---------------------|-------------|
| | C | P | C | P |
| ViP3D (Gu et al., 2023) | 0.226 | - | 2.05 | - |
| GenAD (Zheng et al., 2025) | 0.588 | 0.352 | 0.84 | 0.84 |
| UniAD | 0.495 | 0.361 | <u>0.69</u> | <u>0.79</u> |
| DMAD | <u>0.535</u> | 0.416 | 0.72 | 0.77 |
| SparseDrive | 0.487 | 0.406 | 0.63 | <u>0.73</u> |
| SparseDMAD | 0.500 | <u>0.410</u> | 0.63 | 0.71 |

Table 4: Trajectory prediction results. C and P stand for cars and pedestrians respectively.

the negative transfer in vectorized online mapping, significantly surpassing SparseDrive by 7.0% in mAP, (see Tab. 3).

4.3 Prediction and Planning

Metrics. For motion prediction, we utilize E2E perception accuracy (EPA) proposed in ViP3D (Gu et al., 2023) as the main metric. We also report the minimum average displacement error (minADE). However, since minADE is a true positive metric, it does not fully capture the predictive capabilities of the E2E system, whereas EPA accounts for the number of false positives. For open-loop planning, we use L_2 distances and collision rates. Moreover, we evaluate driving safety in a closed-loop environment using NeuroNCAP (Ljungbergh et al., 2024). This framework reconstructs scenes from the nuScenes dataset and inserts safety-critical objects. The resulting scores are derived from collision rates and impact speeds.

Trajectory prediction. We report car and pedestrian prediction metrics in Tab. 4. Our method surpasses both baselines in EPA, especially achieving improvements of 0.040 for cars and 0.055 for pedestrians over UniAD. However, our method does not improve the minADE of cars. One possible reason is that once detection performance exceeds a certain threshold, further detection improvements often come from reducing false negatives of challenging objects that are either distant or occluded. These hard-to-detect objects typically have limited historical motion data and larger coordinate errors, making them more

| Method | Perception tasks | Ego states in planner | L_2 distances (m) ↓ | | | | Collision rates (%) ↓ | | | |
|--------------------------------|------------------|-----------------------|-----------------------|-------------|-------------|--------------|-----------------------|-------------------|-------------------|--------------------|
| | | | 1s | 2s | 3s | Avg. | 1s | 2s | 3s | Avg. |
| Ego-MLP (Zhai et al., 2023) | ✗ | ✓ | 0.17 | 0.34 | 0.60 | 0.370 | 0 [†] | 0.27 [†] | 0.85 [†] | 0.373 [†] |
| AD-MLP (Li et al., 2024) | ✗ | ✓ | 0.14 | 0.10 | 0.41 | 0.217 | 0.10 | 0.10 | 0.17 | 0.123 |
| VAD (Jiang et al., 2023) | ✓ | ✗ | 0.41 | 0.70 | 1.05 | 0.720 | 0.07 | 0.17 | 0.41 | 0.217 |
| DualVAD (Doll et al., 2024) | ✓ | ✗ | 0.30 | 0.53 | 0.82 | 0.550 | 0.11 | 0.19 | 0.36 | 0.220 |
| GenAD (Zheng et al., 2025) | ✓ | ✗ | 0.28 | 0.49 | 0.78 | 0.517 | 0.08 | 0.14 | 0.34 | 0.187 |
| UniAD* (Hu et al., 2023) | ✓ | ✗ | 0.42 | 0.63 | 0.91 | 0.656 | 0.07 | 0.10 | 0.22 | 0.130 |
| PARA-Drive (Weng et al., 2024) | ✓ | ✗ | 0.25 | 0.46 | 0.74 | 0.483 | 0.14 | 0.23 | 0.39 | 0.253 |
| UniAD | ✓ | ✗ | 0.48 | 0.76 | 1.12 | 0.784 | 0.07 | 0.11 | 0.27 | 0.150 |
| DMAD | ✓ | ✗ | 0.38 | 0.60 | 0.89 | 0.625 | 0.07 | 0.12 | 0.19 | 0.127 |
| SparseDrive | ✓ | ✗ | 0.32 | 0.61 | 1.00 | 0.643 | <u>0.01</u> | 0.06 | 0.22 | <u>0.097</u> |
| SparseDMAD | ✓ | ✗ | 0.30 | 0.61 | 1.01 | 0.643 | 0 | <u>0.07</u> | <u>0.21</u> | 0.093 |

Table 5: **Open-loop planning.** Ego-MLP and AD-MLP are faded since both learn only the ego motion. *Results from the checkpoint in the official repository (UniAD-contributors, 2023), trained with a queue length of 5 in stage 1. †Ego-MLP employs a different strategy in the evaluation of collision rates, therefore the results are not comparable. We reproduce SparseDrive using the official code, but the results differ from its paper because some errors have been fixed after publication.

| Method | NeuroNCAP Scores ↑ | | | | Collision rates (%) ↓ | | | |
|-------------|--------------------|-------------|-------------|-------------|-----------------------|-------------|-------------|-------------|
| | Stat. | Frontal | Side | Avg. | Stat. | Frontal | Side | Avg. |
| UniAD | 3.50 | 1.17 | 1.67 | 2.11 | 32.4 | 77.6 | 71.2 | 60.4 |
| DMAD | 4.40 | 1.47 | 2.07 | 2.65 | 14.8 | 74.0 | 61.6 | 50.1 |
| SparseDrive | 4.42 | 2.96 | 2.30 | 3.23 | 22.4 | 62.8 | 60.4 | 48.5 |
| SparseDMAD | 4.57 | 3.14 | 2.42 | 3.37 | 18.4 | 60.0 | 59.1 | 45.8 |

Table 6: **Closed-loop planning.** We use the official implementation of NeuroNCAP, but our results differ from those in the original paper because the codebase has been updated since its publication.

difficult to predict. A similar issue is observed in UniAD (Hu et al., 2023): in the supplementary materials, UniAD-Large substantially surpasses UniAD-Base in EPA (thanks to better detection and tracking performance), yet it falls short of UniAD-Base in minADE.

Planning. For open-loop evaluation, we adopt the evaluation method of VAD (Jiang et al., 2023), which accommodates the widest range of models to our knowledge. We report our results in Tab. 5. Notably, jointly optimizing L_2 distances and collision rates proves challenging. While PARA-Drive achieves the lowest L_2 distances, it also exhibits the highest collision rates. In the closed-loop evaluation, our structure benefits both baselines in all three cases with stationary, frontal, and side critical objects. We validate that the improvements in perception can be propagated to planning, achieving SOTA collision rates and NeuroNCAP Scores.

4.4 Ablation Study

We ablate our proposed decoders, as shown in Tab. 7, decomposing the motion decoder into three components: motion query, inter-layer, and inter-frame reference point updating.

Model profile. In methods with multi-view camera images as inputs, the primary computational cost is concentrated in the image backbone (Li et al., 2022). In contrast, our approach focuses on the decoding component, resulting in minimal impact on model size and inference speed. Compared to UniAD (Hu et al., 2023), our decoders add 13.1M parameters and increase inference latency by 0.02 seconds on an NVIDIA RTX 6000 Ada.

| Method ID | Interactive semantic dec. | Motion queries | Inter-layer ref. update | Inter-frame ref. update | #Params (M) | Inference time (s) | NDS \uparrow | AMOTA \uparrow | Lanes \uparrow | EPA \uparrow | Avg. $L_2\downarrow$ | Avg. Col. \downarrow |
|-----------|---------------------------|----------------|-------------------------|-------------------------|-------------|--------------------|----------------|------------------|------------------|----------------|----------------------|------------------------|
| 1 (UniAD) | \times | \times | \times | \times | 127.3 | 0.47 | 0.491 | 0.354 | 0.312 | 0.495 | 0.784 | 0.150 |
| 2 | \checkmark | \times | \times | \times | 128.0 | 0.48 | 0.503 | 0.382 | 0.320 | 0.524 | 0.683 | 0.150 |
| 3 | \times | \checkmark | \checkmark | \checkmark | 139.3 | 0.49 | 0.502 | 0.387 | 0.313 | 0.535 | 0.661 | 0.143 |
| 4 | \checkmark | \checkmark | \times | \times | 140.4 | 0.49 | 0.481 | 0.339 | 0.322 | 0.485 | 0.655 | 0.163 |
| 5 | \checkmark | \checkmark | \checkmark | \times | 140.4 | 0.49 | 0.489 | 0.352 | 0.323 | 0.498 | 0.648 | 0.160 |
| 6 | \checkmark | \checkmark | \times | \checkmark | 140.4 | 0.49 | 0.495 | 0.364 | 0.319 | 0.512 | 0.631 | 0.137 |
| 7 (DMAD) | \checkmark | \checkmark | \checkmark | \checkmark | 140.4 | 0.49 | 0.506 | 0.393 | 0.321 | 0.535 | 0.625 | 0.127 |

Table 7: Ablation of DMAD.

Effect of dividing and merging. Experiments ID 1, 2, 3, 7 demonstrate the effectiveness of both proposed decoders. The standalone application of the interactive semantic decoder (ID 2) significantly enhances the performance of object detection, tracking, and map segmentation. The standalone application of the Neural-Bayes motion decoder (ID 3) markedly improves prediction and planning. Notably, ID 3 also significantly enhances detection and tracking, attributed to freeing object queries from learning velocities and the higher-quality reference points provided by the unimodal prediction. Experiments ID 4, 5, 6, 7 show the importance of inter-layer and inter-frame updating in the Neural-Bayes motion decoder.

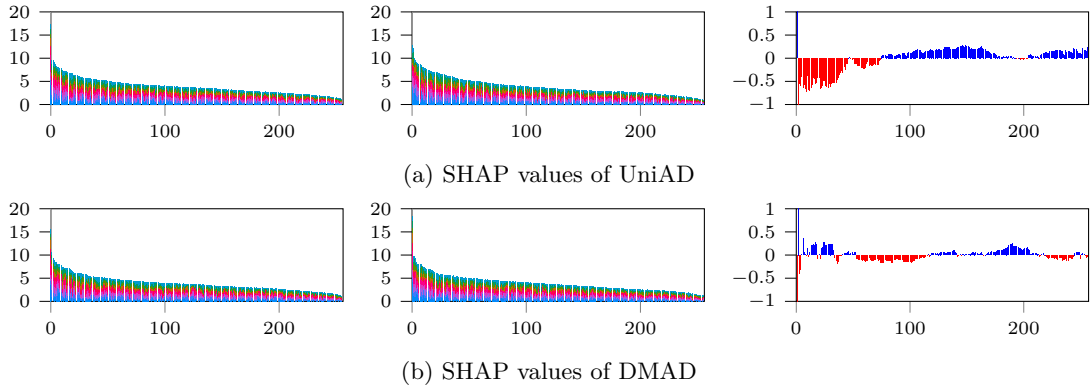
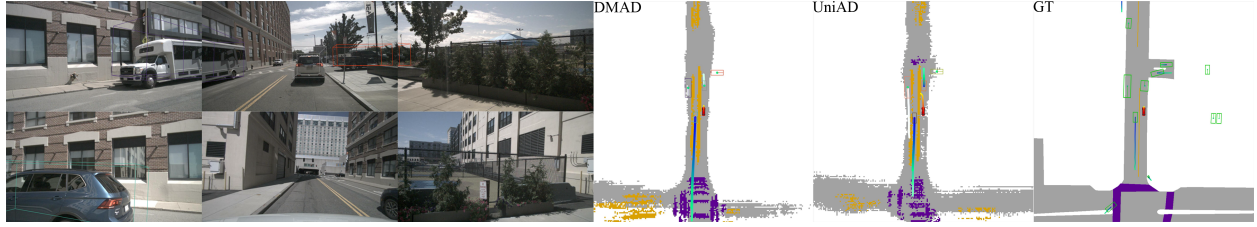


Figure 5: **SHAP values of stage 1 (left), stage 2 (middle), and the difference (right).** Each bar represents the SHAP values of a single feature with respect to different classes. The object query consists of 256 features, forming 256 bars in each chart. The difference is computed as stage 1 minus stage 2, aggregating all classes, where **red** indicates a negative value and **blue** signifies a positive value.

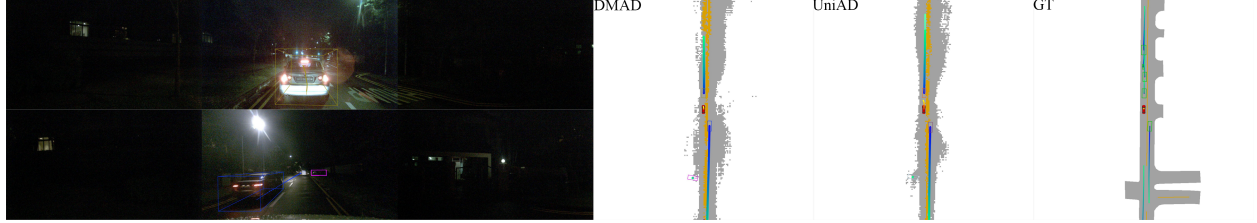
4.5 Visualizations

We use SHAP values (Lundberg & Lee, 2017)—which quantify the contribution of each feature to the change in a model’s output—to inspect the negative transfer in detection and tracking. We visualize the SHAP values of the object query with respect to the object classification output. Changes in SHAP values across the two training stages reveal the negative transfer in UniAD and highlight the effectiveness of our method.

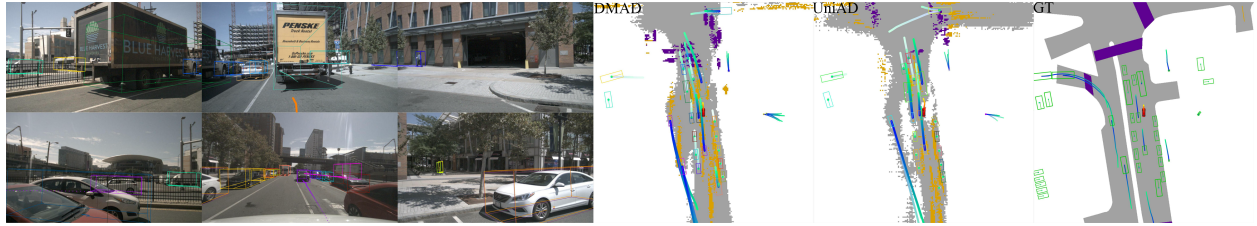
Figure 5a compares the SHAP values between stage 1 and stage 2 of UniAD, sorted in descending order. The left half of the difference bar chart predominantly shows negative values, whereas the right half shows positive values. This indicates that SHAP values in stage 1 are more uniformly distributed, while those in stage 2 are more concentrated. Compared with a flat distribution, this concentration indicates that fewer features are contributing to the classification task, reducing detection and tracking performance. This observation aligns with our argument that during the second stage, object queries are expected to learn motion information, which does not benefit the perception task. Specifically, while the velocity learned in stage 1 is sufficient for tracking (predicting the next timestep), it is inadequate for the long-term prediction over 12 timesteps (6 seconds). Therefore, the object query is forced to learn more motion states that offer limited utility for



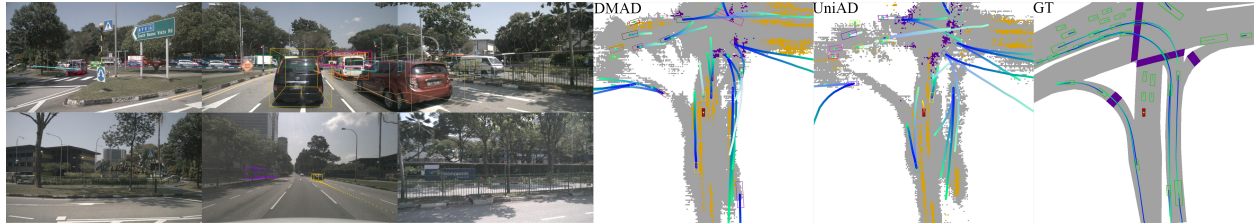
(a) The collision of UniAD is because of an inaccurate prediction of the lead vehicle.



(b) Both models make inaccurate predictions of the lead vehicle during the night. However, UniAD collides with the lead vehicle due to its aggressive driving policy.



(c) An inaccurate detection (the detected position is too close to the ego-vehicle) causes yielding, and then colliding with another vehicle.



(d) UniAD fails to detect the lead vehicle and collides with it.

Figure 6: **Qualitative comparison between DMAD and UniAD.** Each subfigure demonstrates a sample where UniAD encounters collision while DMAD does not.

identifying objects, interfering with the space for semantic information. In contrast, the SHAP values in DMAD maintain a similar distribution across both stages, as shown in Fig. 5b.

We provide qualitative comparisons between DMAD and UniAD in Fig. 6, showcasing how the improved perception and prediction reduces collision rates.

5 Conclusion

In this work, we show that by decoupling semantic and motion learning, we eliminate the negative transfer that E2E training typically imposes on object and map perception. Besides, we leverage the correlation between semantic tasks to promote positive transfer during E2E training. We validate that our improvements in perception and prediction directly enhance planning performance, achieving SOTA collision rates. However, our approach cannot be applied to E2E methods that are without query propagation mechanism, *e.g.*, VAD (Jiang et al., 2023). Addressing this limitation can be our future work.

References

- Mayank Bansal, Alex Krizhevsky, and Abhijit Ogale. ChauffeurNet: Learning to Drive by Imitating the Best and Synthesizing the Worst. In *Robotics: Science and Systems*, 2019.
- Mariusz Bojarski. End to End Learning for Self-Driving Cars. *arXiv preprint arXiv:1604.07316*, 2016.
- Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuScenes: A multimodal dataset for autonomous driving. In *CVPR*, pp. 11621–11631, 2020.
- Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-End Object Detection with Transformers. In *ECCV*, pp. 213–229, 2020.
- Sergio Casas, Wenjie Luo, and Raquel Urtasun. IntentNet: Learning to Predict Intention from Raw Sensor Data. In *CoRL*, pp. 947–956, 2018.
- Sergio Casas, Abbas Sadat, and Raquel Urtasun. MP3: A Unified Model to Map, Perceive, Predict and Plan. In *CVPR*, pp. 14403–14412, 2021.
- Yuning Chai, Benjamin Sapp, Mayank Bansal, and Dragomir Anguelov. MultiPath: Multiple Probabilistic Anchor Trajectory Hypotheses for Behavior Prediction. In *CoRL*, 2019.
- Dian Chen and Philipp Krähenbühl. Learning from All Vehicles. In *CVPR*, pp. 17222–17231, 2022.
- Jiacheng Chen, Yuefan Wu, Jiaqi Tan, Hang Ma, and Yasutaka Furukawa. MapTracker: Tracking with Strided Memory Fusion for Consistent Vector HD Mapping. In *ECCV*, pp. 90–107, 2025.
- Kashyap Chitta, Aditya Prakash, and Andreas Geiger. NEAT: Neural Attention Fields for End-to-End Autonomous Driving. In *ICCV*, pp. 15793–15803, 2021.
- Michael Crawshaw. Multi-Task Learning with Deep Neural Networks: A Survey. *arXiv preprint arXiv:2009.09796*, 2020.
- Henggang Cui, Vladan Radosavljevic, Fang-Chieh Chou, Tsung-Han Lin, Thi Nguyen, Tzu-Kuo Huang, Jeff Schneider, and Nemanja Djuric. Multimodal Trajectory Predictions for Autonomous Driving using Deep Convolutional Networks. In *ICRA*, pp. 2090–2096, 2019.
- Nemanja Djuric, Henggang Cui, Zhaoen Su, Shangxuan Wu, Huahua Wang, Fang-Chieh Chou, Luisa San Martin, Song Feng, Rui Hu, Yang Xu, et al. MultiXNet: Multiclass Multistage Multimodal Motion Prediction. In *IEEE Intelligent Vehicles Symposium (IV)*, pp. 435–442, 2021.
- Simon Doll, Niklas Hanselmann, Lukas Schneider, Richard Schulz, Marius Cordts, Markus Enzweiler, and Hendrik Lensch. DualAD: Disentangling the Dynamic and Static World for End-to-End Driving. In *CVPR*, pp. 14728–14737, 2024.
- Sudeep Fadadu, Shreyash Pandey, Darshan Hegde, Yi Shi, Fang-Chieh Chou, Nemanja Djuric, and Carlos Vallespi-Gonzalez. Multi-View Fusion of Sensor Data for Improved Perception and Prediction in Autonomous Driving. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 2349–2357, 2022.
- Jiyang Gao, Chen Sun, Hang Zhao, Yi Shen, Dragomir Anguelov, Congcong Li, and Cordelia Schmid. VectorNet: Encoding HD Maps and Agent Dynamics From Vectorized Representation. In *CVPR*, pp. 11525–11533, 2020.
- Junru Gu, Chen Sun, and Hang Zhao. DenseTNT: End-to-end Trajectory Prediction from Dense Goal Sets. In *ICCV*, pp. 15303–15312, 2021.
- Junru Gu, Chenxu Hu, Tianyuan Zhang, Xuanyao Chen, Yilun Wang, Yue Wang, and Hang Zhao. ViP3D: End-to-end Visual Trajectory Prediction via 3D Agent Queries. In *CVPR*, pp. 5496–5506, 2023.

- Shengchao Hu, Li Chen, Penghao Wu, Hongyang Li, Junchi Yan, and Dacheng Tao. ST-P3: End-to-end Vision-based Autonomous Driving via Spatial-Temporal Feature Learning. In *ECCV*, pp. 533–549, 2022.
- Yihan Hu, Jiazhi Yang, Li Chen, Keyu Li, Chonghao Sima, Xizhou Zhu, Siqu Chai, Senyao Du, Tianwei Lin, Wenhai Wang, et al. Planning-oriented Autonomous Driving. In *CVPR*, pp. 17853–17862, 2023.
- Bo Jiang, Shaoyu Chen, Qing Xu, Bencheng Liao, Jiajie Chen, Helong Zhou, Qian Zhang, Wenyu Liu, Chang Huang, and Xinggang Wang. VAD: Vectorized Scene Representation for Efficient Autonomous Driving. In *ICCV*, pp. 8340–8350, 2023.
- Zhiqi Li, Wenhai Wang, Hongyang Li, Enze Xie, Chonghao Sima, Tong Lu, Yu Qiao, and Jifeng Dai. BEVFormer: Learning Bird’s-Eye-View Representation from Multi-Camera Images via Spatiotemporal Transformers. In *ECCV*, pp. 1–18, 2022.
- Zhiqi Li, Zhiding Yu, Shiyi Lan, Jiahao Li, Jan Kautz, Tong Lu, and Jose M Alvarez. Is Ego Status All You Need for Open-Loop End-to-End Autonomous Driving? In *CVPR*, pp. 14864–14873, 2024.
- Ming Liang, Bin Yang, Wenyuan Zeng, Yun Chen, Rui Hu, Sergio Casas, and Raquel Urtasun. PnPNet: End-to-End Perception and Prediction with Tracking in the Loop. In *CVPR*, pp. 11553–11562, 2020.
- Bencheng Liao, Shaoyu Chen, Xinggang Wang, Tianheng Cheng, Qian Zhang, Wenyu Liu, and Chang Huang. MapTR: Structured Modeling and Learning for Online Vectorized HD Map Construction. In *ICLR*, 2023.
- Bencheng Liao, Shaoyu Chen, Yunchi Zhang, Bo Jiang, Qian Zhang, Wenyu Liu, Chang Huang, and Xinggang Wang. MapTRv2: An End-to-End Framework for Online Vectorized HD Map Construction. *IJCV*, pp. 1–23, 2024.
- Xuewu Lin, Tianwei Lin, Zixiang Pei, Lichao Huang, and Zhizhong Su. Sparse4D: Multi-view 3D Object Detection with Sparse Spatial-Temporal Fusion. *arXiv preprint arXiv:2211.10581*, 2022.
- Xuewu Lin, Zixiang Pei, Tianwei Lin, Lichao Huang, and Zhizhong Su. Sparse4D v3: Advancing End-to-End 3D Detection and Tracking. *arXiv preprint arXiv:2311.11722*, 2023.
- Yingfei Liu, Tiancai Wang, Xiangyu Zhang, and Jian Sun. PETR: Position Embedding Transformation for Multi-view 3D Object Detection. In *ECCV*, pp. 531–548, 2022.
- Yingfei Liu, Junjie Yan, Fan Jia, Shuailin Li, Aqi Gao, Tiancai Wang, and Xiangyu Zhang. PETRv2: A Unified Framework for 3D Perception from Multi-Camera Images. In *ICCV*, pp. 3262–3272, 2023.
- William Ljungbergh, Adam Tonderski, Joakim Johnander, Holger Caesar, Kalle Åström, Michael Felsberg, and Christoffer Petersson. Neuroncap: Photorealistic closed-loop safety testing for autonomous driving. In *ECCV*, pp. 161–177, 2024.
- Scott M Lundberg and Su-In Lee. A Unified Approach to Interpreting Model Predictions. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *NeurIPS*, pp. 4765–4774. 2017.
- Wenjie Luo, Bin Yang, and Raquel Urtasun. Fast and Furious: Real Time End-to-End 3D Detection, Tracking and Motion Forecasting with a Single Convolutional Net. In *CVPR*, pp. 3569–3577, 2018.
- Tim Meinhardt, Alexander Kirillov, Laura Leal-Taixe, and Christoph Feichtenhofer. TrackFormer: Multi-Object Tracking with Transformers. In *CVPR*, pp. 8844–8854, 2022.
- Jiquan Ngiam, Benjamin Caine, Vijay Vasudevan, Zhengdong Zhang, Hao-Tien Lewis Chiang, Jeffrey Ling, Rebecca Roelofs, Alex Bewley, Chenxi Liu, Ashish Venugopal, et al. Scene Transformer: A unified architecture for predicting future trajectories of multiple agents. In *ICLR*, 2022.
- Chenbin Pan, Burhaneddin Yaman, Senem Velipasalar, and Liu Ren. CLIP-BEVFormer: Enhancing Multi-View Image-Based BEV Detector with Ground Truth Flow. In *CVPR*, pp. 15216–15225, 2024.

- Jonah Philion and Sanja Fidler. Lift, Splat, Shoot: Encoding Images from Arbitrary Camera Rigs by Implicitly Unprojecting to 3D. In *ECCV*, pp. 194–210, 2020.
- Aditya Prakash, Kashyap Chitta, and Andreas Geiger. Multi-Modal Fusion Transformer for End-to-End Autonomous Driving. In *CVPR*, pp. 7077–7087, 2021.
- Shaoshuai Shi, Li Jiang, Dengxin Dai, and Bernt Schiele. Motion Transformer with Global Intention Localization and Local Movement Refinement. *NeurIPS*, 35:6531–6543, 2022.
- Wenchao Sun, Xuewu Lin, Yining Shi, Chuang Zhang, Haoran Wu, and Sifa Zheng. SparseDrive: End-to-End Autonomous Driving via Sparse Scene Representation. *arXiv preprint arXiv:2405.19620*, 2024.
- Sebastian Thrun, Wolfram Burgard, and Dieter Fox. *Probabilistic Robotics*. MIT Press, 2005.
- UniAD-contributors. Planning-oriented Autonomous Driving. <https://github.com/OpenDriveLab/UniAD>, 2023.
- A Vaswani. Attention Is All You Need. In *NeurIPS*, 2017.
- Royden Wagner, Ömer Şahin Taş, Marvin Klemp, Carlos Fernandez, and Christoph Stiller. RedMotion: Motion Prediction via Redundancy Reduction. *Transactions on Machine Learning Research*, 2024.
- Shihao Wang, Yingfei Liu, Tiancai Wang, Ying Li, and Xiangyu Zhang. Exploring Object-Centric Temporal Modeling for Efficient Multi-View 3D Object Detection. In *ICCV*, pp. 3621–3631, 2023.
- Yue Wang, Vitor Campagnolo Guizilini, Tianyuan Zhang, Yilun Wang, Hang Zhao, and Justin Solomon. DETR3D: 3D Object Detection from Multi-view Images via 3D-to-2D Queries. In *CoRL*, pp. 180–191, 2022.
- Xinshuo Weng, Ye Yuan, and Kris Kitani. PTP: Parallelized Tracking and Prediction with Graph Neural Networks and Diversity Sampling. *IEEE Robotics and Automation Letters*, 6(3):4640–4647, 2021.
- Xinshuo Weng, Boris Ivanovic, Yan Wang, Yue Wang, and Marco Pavone. PARA-Drive: Parallelized Architecture for Real-time Autonomous Driving. In *CVPR*, pp. 15449–15458, 2024.
- Chenyu Yang, Yuntao Chen, Hao Tian, Chenxin Tao, Xizhou Zhu, Zhaoxiang Zhang, Gao Huang, Hongyang Li, Yu Qiao, Lewei Lu, et al. BEVFormer v2: Adapting Modern Image Backbones to Bird’s-Eye-View Recognition via Perspective Supervision. In *CVPR*, pp. 17830–17839, 2023.
- Fangao Zeng, Bin Dong, Yuang Zhang, Tiancai Wang, Xiangyu Zhang, and Yichen Wei. MOTR: End-to-End Multiple-Object Tracking with Transformer. In *ECCV*, pp. 659–675, 2022.
- Wenyuan Zeng, Wenjie Luo, Simon Suo, Abbas Sadat, Bin Yang, Sergio Casas, and Raquel Urtasun. End-to-end Interpretable Neural Motion Planner. In *CVPR*, pp. 8660–8669, 2019.
- Jiang-Tian Zhai, Ze Feng, Jinhao Du, Yongqiang Mao, Jiang-Jiang Liu, Zichang Tan, Yifu Zhang, Xiaoqing Ye, and Jingdong Wang. Rethinking the Open-Loop Evaluation of End-to-End Autonomous Driving in nuScenes. *arXiv preprint arXiv:2305.10430*, 2023.
- Lu Zhang, Peiliang Li, Sikang Liu, and Shaojie Shen. SIMPL: A Simple and Efficient Multi-agent Motion Prediction Baseline for Autonomous Driving. *IEEE Robotics and Automation Letters*, 2024.
- Tianyuan Zhang, Xuanyao Chen, Yue Wang, Yilun Wang, and Hang Zhao. MUTR3D: A Multi-camera Tracking Framework via 3D-to-2D Queries. In *CVPR*, pp. 4537–4546, 2022.
- Wenzhao Zheng, Ruiqi Song, Xianda Guo, Chenming Zhang, and Long Chen. GenAD: Generative End-to-End Autonomous Driving. In *ECCV*, pp. 87–104, 2025.
- Zikang Zhou, Luyao Ye, Jianping Wang, Kui Wu, and Kejie Lu. HiVT: Hierarchical Vector Transformer for Multi-Agent Motion Prediction. In *CVPR*, pp. 8823–8833, 2022.

A Tracking as a Semantic Task

We justify the similarity of detection and tracking on nuScenes (Caesar et al., 2020) by analyzing the information learned by the object query. E2E detection and tracking models decode each query into category, location, size, orientation, and velocity. The category is clearly a semantic attribute, while location, size, and orientation serve as spatial complements to the category, all being time-invariant. In contrast, velocity is derived from time, making it a motion attribute. However, measuring velocities is not a common practice in detection, but required by the nuScenes benchmark. Therefore, detection models trained on nuScenes are able to perform tracking without any additional learning effort assuming constant velocity motion (Zhang et al., 2022; Hu et al., 2023; Lin et al., 2023; Gu et al., 2023). Given that current modular E2E models are all trained on nuScenes, we regard the tracking in these methods closely resembles detection, where learning semantics is dominating.

B Bayes Filter

Bayes filter (Thrun et al., 2005) estimates an unknown distribution based on the process model and noisy measurements as follows:

$$p(\mathbf{x}_t | \mathbf{z}_{1:t}) = p(\mathbf{z}_t | \mathbf{x}_t) p(\mathbf{x}_t | \mathbf{z}_{1:t-1}), \quad (3)$$

where \mathbf{x} denotes the state, \mathbf{z} represents the measurement, and the subscript indicates timesteps. The task is to estimate the state \mathbf{x}_t at timestep t given all the measurements $\mathbf{z}_{1:t}$ in the past from timestep 1 to t , which is the product of the likelihood $p(\mathbf{z}_t | \mathbf{x}_t)$ and the prediction $p(\mathbf{x}_t | \mathbf{z}_{1:t-1})$.

Some special cases of Bayes filter, *e.g.*, Kalman filter, are widely used in traditional object tracking. The tracking process can be carried out in three steps: first, predicting the current position based on the object’s historical states $\mathbf{x}_{1:t-1}$; second, identifying the detection most likely to match the prediction as the measurement; finally, updating the current state \mathbf{x}_t according to the latest measurement \mathbf{z}_{t-1} . This process is recursively executed over successive timesteps. We find semantics and motion are similar to the measurement and state in Bayes filter, respectively. Therefore, we introduce the architecture of Bayes filter to transformer decoders, resulting in Neural-Bayes motion decoder.

C Occupancy Prediction

We retain the occupancy prediction module from UniAD to ensure task consistency, where the BEV feature serves as the query and learns from motion prediction features (output queries) through cross-attention. Consequently, we regard occupancy prediction in UniAD as a secondary task to perception and motion prediction, as it merely offers an alternative representation of upstream tasks.

DMAD achieves similar performance (IoU_{near}: 62.7%, IoU_{far}: 39.8%) to UniAD (IoU_{near}: 62.9%, IoU_{far}: 39.6%). The advances of DMAD in upstream tasks do not generalize to occupancy prediction. The reason could be that, by dividing semantics and motion, output features of the prediction module lack spatial information desired by occupancy prediction, such as size, whereas output features of UniAD’s prediction module preserve the spatial information.

D Queue Length

We adopt a different queue length configuration from that of the original UniAD. As mentioned in Sec. 4.1, the rationale behind our decision is that reducing the queue length in stage 2 affects the performance, hindering the observation of negative transfer. Table 8 shows an ablation study of queue length on UniAD, presenting the performance drops by reduced queue length. As the training time scales almost linearly to the queue length, we opt for a queue length of 3 to reduce training time of each iteration.

| Queue length stage 1 | Queue length stage 2 | NDS↑ | mAP↑ | AMOTA↑ | AMOTP↓ | IDS↓ | Lanes↑ | Drivable↑ | EPA↑ | minADE↓ | Avg. L_2 ↓ | Avg. Col.↓ |
|-------------------------|-------------------------|--------------|--------------|--------------|-------------|------------|--------------|--------------|--------------|--------------|--------------|--------------|
| 3 | 3 | 0.491 | 0.377 | 0.354 | 1.34 | 1381 | 0.312 | 0.678 | 0.495 | 0.692 | 0.784 | 0.150 |
| 5 | 3 | 0.499 | 0.381 | 0.362 | 1.34 | 956 | 0.313 | 0.692 | 0.492 | 0.655 | 0.656 | 0.130 |
| 5 | 5 | 0.501 | 0.384 | 0.370 | 1.32 | 885 | 0.314 | 0.690 | 0.495 | 0.714 | 0.615 | 0.123 |

Table 8: **Effect of queue length on UniAD.**

| Unimodal pred. horizon | NDS↑ | mAP↑ | AMOTA↑ | AMOTP↓ | IDS↓ | Lanes↑ | Drivable↑ | EPA↑ | minADE↓ | Avg. L_2 ↓ | Avg. Col.↓ |
|---------------------------|--------------|--------------|--------------|-------------|------------|--------------|--------------|--------------|--------------|--------------|--------------|
| 2s | 0.516 | 0.404 | 0.400 | 1.30 | 695 | 0.321 | 0.691 | 0.534 | 0.735 | 0.679 | 0.220 |
| 4s | 0.506 | 0.396 | 0.393 | 1.30 | 767 | 0.321 | 0.691 | 0.535 | 0.723 | 0.625 | 0.127 |
| 6s | 0.504 | 0.396 | 0.384 | 1.30 | 751 | 0.322 | 0.700 | 0.525 | 0.743 | 0.629 | 0.117 |

Table 9: **Effect of unimodal prediction horizon on DMAD.**

E Effect of Unimodal Prediction Horizon

We conduct experiments on the number of future steps in unimodal prediction, as shown in Tab. 9. We observe that the unimodal prediction horizon influences the proportion of motion information within the BEV feature, thereby impacting the performance of both semantic and motion tasks. A long prediction horizon degrades the performance of semantic tasks, as the BEV feature is forced to prioritize motion learning in order to predict distant future outcomes. Experiments show that a prediction horizon of 6 seconds minimizes the collision rates, but performs worst in tracking. Although this phenomenon can also be referred to as negative transfer, our approach is unable to address this specific type, as the BEV feature is shared across all tasks and is expected to encapsulate both types of information. To balance motion and semantic information within the BEV feature, we set the prediction horizon to 4 seconds.