

S.T.A.R.-Track: Latent Motion Models for End-to-End 3D Object Tracking with Adaptive Spatio-Temporal Appearance Representations

Simon Doll^{1,2}, Niklas Hanselmann^{1,2}, Lukas Schneider¹, Richard Schulz¹,
Markus Enzweiler³, Hendrik P.A. Lensch²

Abstract—Following the tracking-by-attention paradigm, this paper introduces an object-centric, transformer-based framework for tracking in 3D. Traditional model-based tracking approaches incorporate the geometric effect of object- and ego motion between frames with a geometric motion model. Inspired by this, we propose STAR-TRACK which uses a novel *latent motion model* (LMM) to additionally adjust object queries to account for changes in viewing direction and lighting conditions directly in the latent space, while still modeling the geometric motion explicitly. Combined with a novel learnable track embedding that aids in modeling the existence probability of tracks, this results in a generic tracking framework that can be integrated with any query-based detector. Extensive experiments on the nuScenes benchmark demonstrate the benefits of our approach, showing *state-of-the-art* (SOTA) performance for DETR3D-based trackers while drastically reducing the number of identity switches of tracks at the same time.

Index Terms—Visual Tracking, Deep Learning for Visual Perception, Autonomous Vehicle Navigation

I. INTRODUCTION

ROBUST perception and tracking of movable objects in the environment form the basis for safe decision-making in autonomous agents such as self-driving cars. Classical *multi-object tracking* (MOT) pipelines follow a *tracking-by-detection* paradigm, using object detectors coupled with greedy matching [1] and state estimators [2], [3] to track objects. Building on recent advances in object detection from multi-view camera images, transformer-based architectures [4], [5], [6] can yield strong tracking performance [1], [7] using relatively low-cost sensors. However, decoupling the detection and tracking tasks comes with two main drawbacks: (1) the object detection model is optimized towards a detection metric, rather than directly optimizing for the downstream tracking performance, which is prone to compounding errors [8], [9] and (2) it makes it non-trivial to incorporate appearance information, which poses a challenge to consistent association. This in particular can lead to difficulties in handling confusion

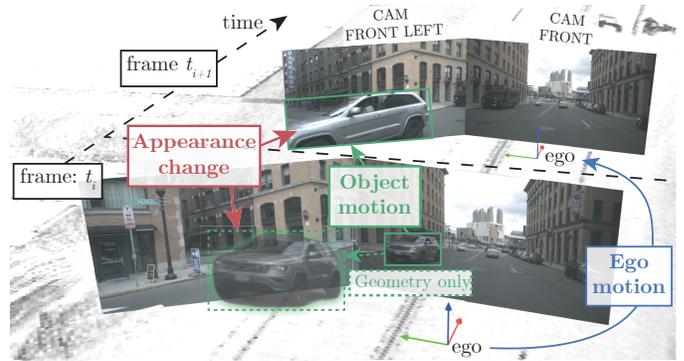


Fig. 1. Visualization of a tracked object for two consecutive frames. Due to ego and object motion the 3D pose and the appearance of the object in the camera images change in scale, viewing angle and lighting condition. We utilize an explicit geometric and a novel latent motion model to compensate for these effects during the prediction step of the tracking pipeline.

among object identities in crowded scenarios with many partial object-to-object occlusions [10].

Recent works [11], [12] propose an alternative *tracking-by-attention* paradigm that unifies perception and tracking into a single module. Under this paradigm, the geometric and semantic information contained in the latent object queries of query-based detectors can be leveraged for the association of object instances across time via attention [13]. Additionally, tracking-by-attention allows to use these queries as detection priors in the following frames. This requires adjusting them to the expected future object state, analogous to the model-based prediction step in classical state estimator-based trackers [14].

For geometric features, this can be done by applying the transformation corresponding to both ego and estimated object motion. However, this is not possible for latent object queries, as they also encode semantics and appearance in addition to geometric information. MUTR3D [12] sidesteps this issue by anchoring object queries to geometric reference points which can be analytically updated. While this enables some adjustment, only the object translation rather than the full pose is considered and the change in appearance resulting from changes in the relative pose is not modeled. A tracking method that corrects both geometric and appearance information directly in latent space is proposed in [15]. However, this approach forfeits the ability to analytically update geometric information and does not model object motion.

Manuscript received: 7, 26, 2023; Revised 10, 9, 2023; Accepted 12, 5, 2023.

This paper was recommended for publication by Editor Ashis Banerjee upon evaluation of the Associate Editor and Reviewers' comments. This work was supported by the German Federal Ministry for Economic Affairs and Climate Action (KI Delta Learning: Förderkennzeichen 19A19013A).

¹Mercedes-Benz AG simon.doll@mercedes-benz.com

²University of Tübingen

³Institute for Intelligent Systems, Esslingen University of Applied Sciences
Digital Object Identifier (DOI): 10.1109/LRA.2023.3342552

In this paper, we propose to compensate for appearance changes resulting from both ego and object motion via a novel LMM which updates queries in latent space as a function of the geometric motion. Paired with analytical updates on geometric reference points for each query, we obtain a transformable **Spatio-Temporal geometry and Appearance Representation** for each object that enhances consistency with future observations. Furthermore, we propose track embeddings that encode information on the lifetime of a tracked object to distinguish track queries from new detections. Our approach termed STAR-TRACK, exhibits improved tracking performance. Specifically, we observe that accounting for appearance changes between frames as well as the improved existence probability modeling eases association, leading to a drastically reduced number of switches in object instance identities.

In summary, we make the following main contributions:

- We are the first to compensate for the appearance change induced by both, ego- and object motion in a tracking-by-attention paradigm leveraging a *latent motion model* (LMM) that extends query-based object detectors.
- We introduce novel track embeddings allowing to implicitly model the life cycle of a tracked object.
- We outperform current state of the art DETR3D-based tracking approaches on nuScenes [16] where the LMM and track embeddings in particular reduce fragmentations and identity switches by a large margin.

II. RELATED WORK

Query-based Detection: MOT approaches that follow the tracking-by-detection [1], [2], [3] paradigm require a detector to detect a set of objects in each frame. The pioneering work DETR [17] proposed a way to leverage the transformer architecture for object detection. In contrast to previous approaches, this set-based architecture comes with various desirable properties such as a sparse prediction scheme, a dynamic amount of object hypotheses, and no need for hand-crafted components such as *non-maximum suppression* (NMS). Additionally, the concept was generalized to the 3D case as well as to different sensor modalities including LiDAR [18], [19], multi-view camera [4], [5] and multi-modal detection methods [1], [18]. It is noteworthy that such query-based detectors became the de-facto standard in object detection and reach SOTA performance on various benchmarks such as COCO [20] or nuScenes [16].

Tracking-by-Detection: Methods that rely on the well-established tracking-by-detection paradigm have the benefit of being compatible with any detection framework since the detection per frame and the tracking/association part are not directly linked. A simple greedy association [21] is still widely adopted in SOTA methods on the nuScenes tracking benchmark [1], [7], [18]. In this generic approach, the detector can not make use of previous tracks and the association often relies on geometric cues only. This causes track identity switches in which a track is reinitialized with a detection instead of the detection being associated with the previous track. Various extensions such as re-ID features [22], [23]

and motion models [3], [24] have been proposed to mitigate this effect. Motion models integrate prior knowledge about the physical properties and trajectory of the tracked object while re-ID features allow an association that is not solely based on bounding box geometry but also influenced by other features such as motion cues or objects appearance.

Tracking-by-Attention: To overcome the independent nature of the detection and tracking modules in a fully differentiable fashion and to implicitly solve the association between frames, the *tracking-by-attention* paradigm can be used [10], [25]. Leveraging the potential of attention, tracking and detection are performed jointly by auto-regressive query-based tracking since each detection of the last frame is used as a prior (track-query) for the next frame. MUTR3D[12] extends the object detection method DETR3D [4] for tracking by adding a geometric compensation of object and ego motion. This is done by utilizing a 3D reference point per object that is transformed between consecutive frames while the latent query features remain unchanged. A possibility to account for an appearance change caused by the ego-motion is presented in [15]. The proposed ego-motion-compensation module models the effect directly in the latent space as a linear function that depends on the estimated transformation between the two time steps. Similar to the 3D case in which the transformation can be represented as a homogeneous matrix this transformation in latent space is modeled as a full-rank matrix which is learned from the given ego motion via a hyper-network [26].

Inspired by the aforementioned previous works, we propose a latent motion model to account for the effects of *ego and object motion* on the latent appearance representation jointly. This allows for keeping the explicit geometric update proposed in [12] while altering the learned appearance of an object as a function of the geometric transformation to simplify its detection and re-identification in the next frame.

III. METHOD

Our proposed approach tackles *multi-object tracking* from multi-view camera images. Given a set of c mono camera images $\mathbf{I}_c \in \mathcal{I}_t$ with shape $H \times W \times 3$ for each timestamp t the tracking objective is to estimate a set of bounding boxes $\mathbf{b}_t^{id} \in \mathcal{B}_t$ with $\mathbf{b}_t^{id} = [x, y, z, w, l, h, \theta, v_x, v_y]$ describing each object as defined in nuScenes [7]. Besides center, shape, heading angle and velocity of each object, each bounding box has a corresponding *id* that is consistent over time.

A. Overall Architecture

Under the tracking-by-attention paradigm, an object is tracked by updating its unique query feature to be consistent with new observations and other track hypotheses at each point in time via the attention mechanism [10], [12]. Since attention reasons about the affinity between new observations and existing tracks via feature similarity, queries of tracked objects need to be adjusted to account for the changes in relative pose and appearance resulting from both ego- and object motion between frames. To this end, we propose an LMM, an extension to commonly used purely geometric motion

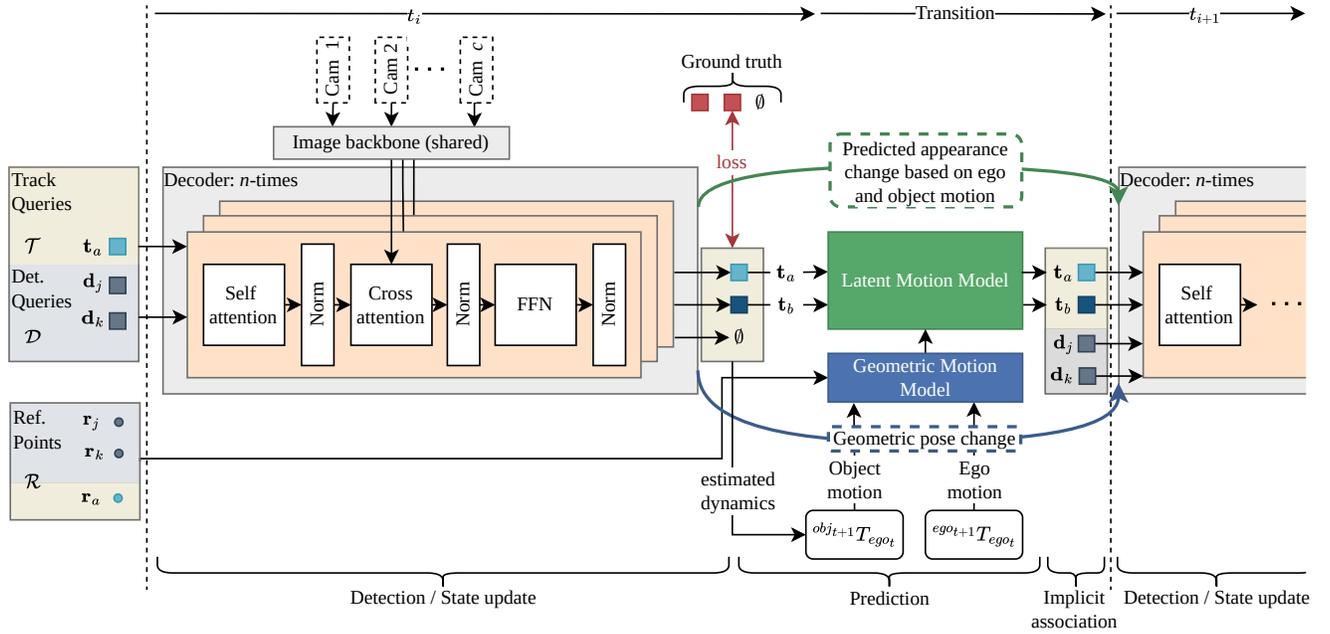


Fig. 2. STAR-TRACK architecture. A joint set of time-independent object queries and track queries of the previous frames is used in a stack of decoder layers that utilize self- and cross-attention blocks to detect and re-identify objects in consecutive time steps. This requires predicting the state of each object in the following frame. Combined with any geometric motion model (blue) the newly proposed latent motion model (green) solves this issue by modeling the spatio-temporal change of a track query in the latent and the 3D geometric space jointly, based on the estimated dynamics.

models. The LMM adjusts the latent feature of each object to be consistent with the expected state in the next frame, increasing similarity to new observations of the same object and simplifying the association task. The LMM implements a generic query prediction strategy that can be readily coupled and jointly trained with any query-based detector.

An overview of the proposed architecture is presented in Fig. 2. We utilize a decoder-only transformer architecture as in DETR3D [4], where a set of learnable detection queries $\mathcal{D} = \{d_1, \dots, d_n\}$ is used to represent hypotheses for newly detected objects in the scene. Following the MUTR3D [12], the time independent detection queries are concatenated with a set of track queries $\mathcal{T} = \{t_1, \dots, t_m\}$ that correspond to hypotheses from the previous frames. Then, the decoder refines both the track hypotheses and new detections jointly by applying self- and cross-attention into features extracted from multi-view camera images \mathcal{I}_t by a shared image backbone in an alternating fashion. We kindly refer the reader to [11], [12], [17] for further details on the general MOT architecture. Lastly, the bounding boxes \mathcal{B}_t are obtained with a *feed forward network* (FFN) while we carry the objects over to the next frame by applying both the analytical geometric motion transformation as well as the LMM.

B. Revisiting Multi-Object Tracking

Model-based tracking systems [2], [27] typically rely on sequential steps that allow to incorporate inductive biases into the different parts of the tracking framework while also maintaining a high level of interpretability.

Detection / State update: In each frame a set \mathcal{D} of new detections is used to update the current belief state of tracked

objects \mathcal{T} in the scene. This enables rejecting implausible sensor measurements, updating the estimated bounding box and existence probability of each track, and spawning new tracks for newly appeared objects. The transformer-based tracking-by-attention mechanism mirrors this behavior by performing two attention operations utilizing *scaled dot product attention* [13]. *Self-attention* within the joint set of track queries and newly spawned detection queries models object interactions, integrating new objects and rejecting duplicate proposals. Subsequently, *cross-attention* between all object queries and the camera features is used to refine each object proposal by incorporating sensor measurements. The tracking-by-attention framework utilizes track queries of previous time steps as priors for the detection in the next frame which potentially simplifies the detection of objects that are far away, partially occluded, or hardly visible.

Prediction: Given the current ego motion transformation ${}^{ego_{t+1}}T_{ego_t}$ Eq. (1) and estimated object dynamics, for instance the velocity of each tracked object, a traditional geometric tracking framework predicts the object pose in the next frame. This is typically achieved utilizing a motion model which is a function of object state and dynamics.

For a latent object representation the geometric update in terms of the object pose should be handled similarly to the explicit bounding box representation since the geometric transformation can be applied analytically. However, the high-dimensional appearance representation of the object query also needs to be taken into consideration since the ego and object motion might heavily affect the appearance of the object and thus its query feature in the next frame, see Fig. 1. This is crucial since the transformer attention relies on a query-

key similarity [13]. Without a latent appearance update the re-identification of a tracked object in the next frame might be impaired. Firstly, track identity switches or track losses can occur if a track query cannot be associated to the sensor data of the next frame in the cross-attention blocks. Secondly, without proper appearance updates, duplicates might spawn, since existing tracks fail to suppress their newly detected counterparts in the self-attention blocks.

Association: To associate detections in the next frame with existing tracks, any similarity metric between object hypotheses can be used. Traditional methods rely on geometry-based metrics [21], [28] or additional re-ID features [22], [29] to form an affinity matrix between tracks and new detections which can be used together with the Hungarian algorithm to find an optimal matching. Auto-regressive query-based tracking methods [10], [11], [12], [25] solve this problem differently since a track query always represents the same object in the scene resulting in an implicit association. During training, this is enforced by matching each track query to its corresponding object in the scene to which it was assigned at first appearance. If two hypotheses describe the same object, the model needs to distinguish between newly spawned and already tracked objects and favor the latter. This is crucial since confusions between tracks and newborn detections might result in track losses or identity switches between tracks and new detections at inference time.

As a result of the considerations above, two key challenges arise for auto-regressive query-based tracking: (1) The prediction step needs to model the influence of the geometric transformation on the pose of the object as well as its latent appearance and semantic features. (2) Due to the implicit association mechanism each track query needs a latent existence probability to efficiently suppress newborn duplicate queries that also belong to the tracked object.

C. Latent Motion Models

The prediction step in the tracking pipeline aims to estimate the state of an object in the next timestep. In our model, the set of tracked objects and newly spawned detections is defined as a set of latent vectors $\mathbf{q} \in \mathcal{Q}$. Additionally, the position of each object query is defined with respect to a 3D reference point $\mathbf{r} \in \mathcal{R}$ as proposed in [4]. Consequentially, the geometric effect of the ego motion for a time delta δ_t between two frames can be described with a homogeneous matrix that combines rotation \mathbf{R} and translation \mathbf{t}

$${}^{ego}_{o_{t+1}}\mathbf{T}_{ego_t} = \begin{bmatrix} \mathbf{R} & \mathbf{t} \\ 0 & 1 \end{bmatrix}. \quad (1)$$

Furthermore, the regression branches of the transformer decoder predict an estimate of the dynamics for each object. These include the estimated velocity $\mathbf{v} = (v_x \ v_y)$, that is supervised by ground truth data during training, and an optional turn-rate δ_θ for the heading angle θ resulting in

$${}^{e'_t}\mathbf{T}_{e_t} = \begin{bmatrix} \cos(\delta_\theta) & -\sin(\delta_\theta) & 0 & v_x \cdot \delta_t \\ \sin(\delta_\theta) & \cos(\delta_\theta) & 0 & v_y \cdot \delta_t \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}. \quad (2)$$

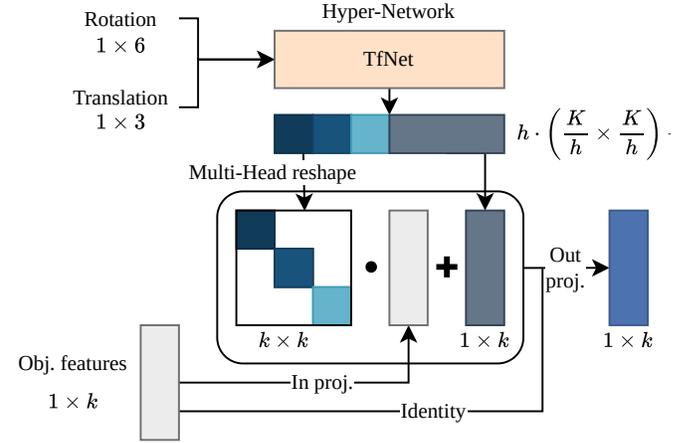


Fig. 3. Latent motion model architecture. A geometric transformation consisting of a translation and rotation is applied to the high-dimensional object query by using a sparse latent transformation matrix K . We estimate the elements of K with a hyper-network (TfNet) and apply the transformation as an input dependent multiplication, mimicking the behavior of a homogeneous matrix in 3D. Note the sparse block-diagonal shape of the generated matrix.

For consistent notation, we propose an auxiliary frame e'_t that describes the state of the world after object motion compensation relative to the ego frame at time t . We note that due to the explicit modeling of this transformation, any motion model [28] can be used to constrain the estimated transformations by model-based assumptions.

Hyper-Networks: Besides the explicit geometric update on the reference point \mathbf{r} of an object as defined in Eq. (3), an additional update to the latent features \mathbf{q} is required to propagate those to the next frame. As argued in [15], the effect of the geometric transformation in latent space can be modeled as a linear operator that performs an input-dependent multiplication on the object query in the form of a latent transformation matrix ${}^b\mathbf{K}_a$. This matrix is a function of its geometric counterpart ${}^b\mathbf{T}_a$ and represents an arbitrary transformation from frame a to frame b . Geometric and latent information is jointly updated:

$$\mathbf{r}_{e_{t+1}} = {}^{e_{t+1}}\mathbf{T}_{e'_t} \cdot {}^{e'_t}\mathbf{T}_{e_t} \cdot \mathbf{r}_{e_t} \quad \text{Geometric Update} \quad (3)$$

$$\mathbf{q}_{e_{t+1}} = {}^{e_{t+1}}\mathbf{K}_{e'_t} \cdot {}^{e'_t}\mathbf{K}_{e_t} \cdot \mathbf{q}_{e_t} \quad \text{Latent Update} \quad (4)$$

We propose a *transformation hyper-network* (TfNet) to estimate the parameters of the latent transformation matrix ${}^b\mathbf{K}_a$. This matrix is applied as an input-dependent multiplication with the latent object query \mathbf{q} . A latent translational offset is incorporated as an element-wise addition. An overview of the proposed LMM architecture is given in Fig. 3.

Input Representation: The input to the TfNet consists of a rotational and translational part:

$${}^b\mathbf{K}_a = \text{TfNet}({}^b\mathbf{R}_a, {}^b\mathbf{t}_a), \quad (5)$$

whereas ${}^a\mathbf{R}_b$ describes the rotational component and ${}^a\mathbf{t}_b$ the translation of ${}^a\mathbf{T}_b$. While the translation is represented as a 3D vector, we utilize the 6D rotation representation proposed in [30] to increase the numeric stability.

Sparse Latent Transforms: Since the latent features are typically high-dimensional [4], [12], a hyper-network that predicts ${}^b\mathbf{K}_a$ as a full-rank matrix might be over-parameterized or even intractable to train. This is due to the large number of parameters in the output weight matrix ${}^b\mathbf{K}_a$ that need to be computed per object in each frame. We mitigate this potential issue by adopting the concept of multi-head attention from [4], [13], [17] and propose a sparse multi-head LMM. Here, attention is computed as a combination of h different low-dimensional attention heads that operate on h splits of the feature vector with a dimensionality of $h_{dim} = k/h$ each. Instead of predicting $k^2 = h^2 \cdot h_{dim}^2$ weights for a full-rank description of K , we propose to only predict $h \cdot h_{dim}^2$ weights for a sparse approximation that drastically reduces the parameter count of the latent transformation matrix. Analogously to the attention computation, these are then used as heads along the diagonal of ${}^b\mathbf{K}_a$ that operate on parts of the k -dimensional latent vector \mathbf{q} , see Fig. 3. Since only neighboring dimensions of the feature vector that lie within the same head can influence the latent transform, we follow [13] and incorporate an input and output projection to mitigate this effect.

As a result, with the multi-head LMM the latent transformation can be directly applied to the full latent vector in a sparse and numerically more stable fashion, while also streamlining the architecture to follow the layout of the attention blocks that are used in all other parts of the model.

D. Track Embeddings for implicit Existence Probability

As discussed in Section III-B, the self-attention blocks serve the purpose of allowing for object interactions as well as suppressing newborn detections that belong to an already tracked object. Although it might be sufficient to distinguish between tracks and new detections in this case, the track queries in general require a consistent integration of the track history to account for short-term occlusions and deliver robust existence probability estimates.

Since learned embeddings have been used successfully to incorporate inductive biases in attention-based detectors [18], [5], we propose to use a learned latent *track embedding* to address the aforementioned issues. Using a single shared track embedding \mathbf{e} and a FFN we update all active tracks \mathcal{T} of the current time step using

$$\mathbf{t}_i' = \mathbf{t}_i + \text{FFN}([\mathbf{t}_i, \mathbf{e}]) \quad \forall \mathbf{t}_i \in \mathcal{T}. \quad (6)$$

This way, the model is flexible to integrate the track embedding to the current latent state of an object and to model the desired distinction between tracks and new detections. As a result, we obtain more consistent existence probabilities and improved track losses, track fragmentations and identity switches, as our experiments in Section IV-B show.

IV. EXPERIMENTS

We evaluate the performance of STAR-TRACK on the tracking task [7] of the nuScenes dataset [16]. Additionally, we provide extensive ablation studies to evaluate the effects of different LMM configurations, latent track embeddings and transform representations, as well as qualitative results.

A. Experimental Setup

Dataset: All experiments are performed on the large-scale nuScenes dataset [16] that consists of 1000 scenes with a length of 20s with a frequency of 2Hz. We use the official training, validation and test set split and the seven object classes required for the tracking benchmark [7].

Metrics: We report performance using the metrics as defined in the nuScenes benchmark [7]: These include the *average multi object tracking accuracy* (AMOTA) as well as the *average multi object tracking precision* (AMOTP). Additionally, we report the *number of identity switches* (IDS), *number of track fragmentations* (FRAG) and *number of mostly tracked trajectories* (MT) as secondary metrics. For the full metric definitions and further details, we refer to [7], [16].

Training Configuration: To increase comparability and reproducibility, we closely follow the settings proposed in MUTR3D [12]. Each training sample consists of three consecutive frames. The geometric and latent motion models assume a constant velocity and no turn-rate transformation for each object, as used in [12]. We leave the integration of more complex dynamics models to future work. As in previous works [4], [12], bi-partite matching and the Hungarian algorithm are used to match tracked objects of the current frame with the ground truth. We use Focal-Loss [31] as classification loss and L1-Loss for bounding box regression. In the training phase, previously matched track queries are always matched to their corresponding ground truth objects. As in [11], [12], we drop tracked queries with a probability $p_{drop} = 0.1$ and spawn false positive tracks with a probability of $p_{fp} = 0.3$. During inference, non-confirmed tracks are kept as inactive for five frames to handle full occlusions over multiple time steps.

We train all models for 24 epochs with the same random seed on four NVIDIA-V100 GPUs using a batch size of four and a ResNet-101 backbone [32]. As proposed in [12], the transformer utilizes $l = 6$ decoder layers, $q = 300$ detection queries for each frame and a latent dimension of $d_l = 256$ spread over $h = 8$ heads of dimension $d_h = d_l/h = 32$. This is also used as configuration of the proposed LMM. All experiments use the training schedule proposed in DETR3D [4] that utilizes a learning rate of $2e^{-4}$, a cosine annealing learning rate schedule and AdamW [33].

We initialize the model with an already trained MUTR3D checkpoint to avoid retraining and keep the image backbone and *feature pyramid network* (FPN) fixed. To initialize the newly introduced LMM, we propose a simple yet effective pretraining scheme: For each sample in the dataset we store the tracking results, consisting of latent queries as well as decoded object proposals from MUTR3D [12] and train the LMM to predict the state of the latent object query vectors of the next frame.

B. Comparison to Existing Works

We compare STAR-TRACK to state-of-the-art methods for 3D MOT on multi-view camera images. To control for the effects of different detection algorithms on the overall tracking

TABLE I

COMPARISON OF STATE-OF-THE-ART METHODS ON THE nuSCENES BENCHMARK. FOR A FAIR COMPARISON ALL METHODS ON THE VALIDATION SET UTILIZE DETR3D [4] AS DETECTOR. DETR3D[†] UTILIZES THE GREEDY TRACKING APPROACH PROPOSED IN [21], DETR3D[‡] THE MORE ELABORATE TRACKING APPROACH INTRODUCED IN [34]. DUE TO A POTENTIAL EVALUATION ERROR IN MUTR3D [12], [35] WE ADD A CUSTOMIZED MUTR3D⁺ BASELINE. THE VERSION OF OUR MODEL THAT ONLY USES THE LMM AND NO LEARNED TRACK EMBEDDING IS DENOTED BY *.

Name	Backbone	#Params	AMOTA [↑]	AMOTP [↓]	RECALL [↑]	MOTA [↑]	MT [↑]	FRAG [↓]	IDS [↓]	FPS [↑]
Validation-Split										
DETR3D [4] [†]	ResNet101	-	0.327	1.372	0.463	0.291	2039	2372	2712	-
DETR3D [4] [‡]	ResNet101	-	0.353	1.382	0.469	0.315	2065	2309	1807	-
MUTR3D [12]	ResNet101	59M	0.294	1.498	0.427	0.267	-	-	3822	6.98
MUTR3D [12] ⁺	ResNet101	59M	0.360	1.411	0.487	0.341	2368	1232	522	6.98
CC-3DT [36]	ResNet101	-	0.359	1.361	0.498	0.326	-	-	2152	2.06
PF-Track [37]	VovNet-V2-99	-	0.362	1.363	-	-	-	-	300	-
STAR-TRACK*	ResNet101	62M	0.378	1.365	0.497	0.354	2467	1241	439	6.86
STAR-TRACK	ResNet101	62M	0.379	1.358	0.501	0.360	2468	1109	372	6.76
Test-Split										
MUTR3D [12]	ResNet101	59M	0.270	1.494	0.411	0.245	2221	2749	6018	-
STAR-TRACK	VovNet-V2-99	83M	0.439	1.256	0.562	0.406	3726	1250	607	6.68

performance, we present our main comparison in terms of DETR3D-based frameworks, which are well-established and widely used [5], [38], [6]. This allows for a fair assessment of our contributions.

As shown in Table I, our tracking framework STAR-TRACK that utilizes the novel LMM and track embedding achieves the best performance in all key metrics on the nuScenes benchmark [7] for DETR3D-based [4] tracking algorithms without reducing the inference speed.

In comparison to the greedy tracking DETR3D baseline that uses a purely geometry-based prediction and association [21], our framework improves the main metric AMOTA substantially by 5.2%. The optimized version of MUTR3D [12], [35] is outperformed by 1.9%, highlighting the crucial role of the LMM. In particular, we observe a drastic reduction of IDS by 86.2% compared to the greedy version and by 28.7% compared to MUTR3D, see Table I. We address this fact to the spatially and temporally consistent appearance representations provided by the LMM and our proposed track embedding. This benefits the association resulting in less track fragmentations (FRAG) and a higher amount of mostly tracked trajectories (MT).

Additionally, STAR-TRACK also outperforms the concurrent works PF-Track [37] by 1.7% and CC-3DT [36] by 2% AMOTA, respectively. The former employs advanced query refinement operations for temporal consistency and a stronger VovNet-V2-99 [39] image backbone, the latter proposes a LSTM-based learned motion model [36]. Evaluating our model with a VovNet-V2-99 trained on both the train and validation set on the nuScenes test set results in 43.9% AMOTA. This improves over MUTR3D [12] by 16.9% and even outperforms concurrent work that utilizes stronger detection algorithms [36], [37].

C. Ablation and Analysis

Qualitative results: A qualitative example of two consecutive time steps is shown in Fig. 4. STAR-TRACK is particularly strong in handling large appearance changes, e.g. due to different lighting conditions and tracking road participants

TABLE II

EFFECT OF TRAINING TIME. FOR A FAIR COMPARISON WE FINE-TUNE OUR VERSION OF MUTR3D [12] WITH AND WITHOUT AN LMM INDICATED BY W/LMM. RUNS DENOTED BY W/INIT USE A PRETRAINED MUTR3D INSTEAD OF A PRETRAINED DETR3D [4] CHECKPOINT.

w/LMM	w/Init	AMOTA [↑]	AMOTP [↓]	IDS [↓]
X	X	0.338	1.425	531
X	✓	0.358	1.382	492
✓	✓	0.378	1.365	439

TABLE III

EFFECT OF DIFFERENT LMM ARCHITECTURES. W/LMM INDICATES WHETHER AN LMM IS USED, MULTI-HEAD (W/MH) DENOTES A SPARSE LATENT TRANSFORMATION MATRIX ${}^b\mathbf{K}_a$ INSTEAD OF A FULL-RANK VERSION. THE HEAD / MATRIX SIZE IS DENOTED BY $|K|$.

w/LMM	w/MH	$ K $	AMOTA [↑]	IDS [↓]
X	-	-	0.358	492
✓	X	32 ²	0.372	432
✓	X	96 ²	0.370	402
✓	✓	16 · 16 ²	0.374	517
✓	✓	4 · 64 ²	0.373	434
✓	✓	8 · 32 ²	0.378	439

under strong object-object occlusions. We provide additional videos of the tracking performance in the supplementary.

Effect of Training Time: The effect of longer training schedules is shown in Table II. MUTR3D [12] gains a performance boost of 2.0% in AMOTA and 7.3% in IDS by further fine-tuning. Adding the proposed LMM yields 2% AMOTA and improves the IDS by 10.7% as compared to the equally long trained model. This clearly indicates that the use of our LMM results in more consistent tracks with a reduced number of identity switches.

Effect of LMM Architecture: The performance of different LMM architectures is shown in Table III. Using the proposed sparse multi-head LMM instead of a full-rank representation of the latent motion matrix ${}^b\mathbf{K}_a$ does not only align the architecture to the multi-head attention blocks but also reduces the amount of output parameters of the hyper-network. This

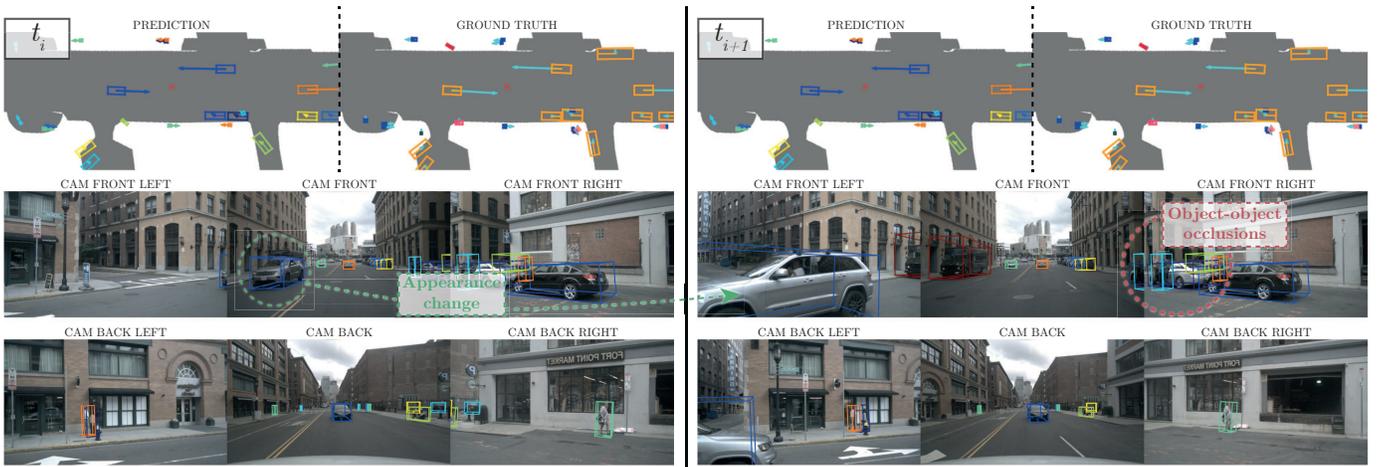


Fig. 4. Qualitative results for two consecutive frames on the nuScenes [16] validation set. Upper row shows predictions and ground truth in top view. Different colors of the predicted objects indicate different object ids. The bottom row shows the predictions projected to the multi-view camera images.

is key to scale the latent transformation matrix to the full latent space dimensions. Using the same configuration as the attention blocks of the transformer for the multi-head LMM results in an boost in AMOTA of 0.8% over a LMM that uses a full-rank latent motion matrix.

Effect of Transform Representation: Different strategies to apply the transformation modeled by the LMM are shown in Table IV. We do not observe a performance increase when the latent query is used as an additional input to the TfNet. This is in line with our general design paradigm to compute the latent motion matrix solely from its geometric counterpart. Although it is beneficial to apply the LMM twice instead of merging object and ego motion, using shared parameters for the object and ego motion compensation cuts the number of parameters in half and does not cause any ill-effects. This supports the general design to model any geometric transformation with the LMM without creating an explicit distinction between object and ego motion.

Integration to other methods: To showcase the flexibility of the proposed LMM we incorporate it into the concurrent work StreamPETR [40] that proposes motion-aware layer normalization (MLN). As shown in Table V our proposed architecture improves the NDS by 1.3% when using the LMM which is a generalized version of the MLN.

Inference latency: An analysis of the runtime of different components of STAR-TRACK is shown in Table VI. The proposed LMM only adds additional 1.48% latency and the track-embeddings 1.52% respectively, since the runtime is dominated by the image backbone and transformer layers for both DETR3D [4] and StreamPETR-based [40] models.

V. CONCLUSION

This paper presented STAR-TRACK, a novel approach for 3D object tracking-by-attention that is compatible with any query-based object detector. We transferred the concept of motion models from traditional geometry-based trackers to the tracking-by-attention paradigm in terms of latent motion

TABLE IV
LMM TRANSFORM REPRESENTATION. MODELS THAT APPLY OBJECT AND EGO MOTION SEPARATELY ARE DENOTED WITH W/SEPARATE. W/SHARE INDICATES MODELS THAT USE SHARED PARAMETERS AND W/FEATS LMMs THAT UTILIZE THE QUERY FEATURE AS INPUT TO THE TfNET.

w/Separate	w/Share	w/Feats	AMOTA \uparrow	IDS \downarrow
\times	\checkmark	\times	0.370	492
\times	\checkmark	\checkmark	0.371	446
\checkmark	\times	\times	0.377	411
\checkmark	\times	\checkmark	0.366	464
\checkmark	\checkmark	\checkmark	0.370	426
\checkmark	\checkmark	\times	0.378	439

TABLE V
PERFORMANCE OF STREAMPETR [40] ON THE VALIDATION SET OF THE NUSCENES DETECTION BENCHMARK. \clubsuit INDICATES A MODEL THAT USES THE PROPOSED LMM INSTEAD OF THE MLN.

Name	mAP \uparrow	mATE \downarrow	mAOE \downarrow	mAVE \downarrow	NDS \uparrow
StreamPETR	0.483	0.591	0.479	0.195	0.562
StreamPETR \clubsuit	0.485	0.611	0.367	0.185	0.575

models that predict the spatio-temporal appearance change of objects between two frames. This allowed for a prediction step that models a geometric transformation in an analytical way and applies this transformation in the latent space with a learned motion matrix at the same time. An additional latent track embedding improved the latent existence probability of tracks. In our experimental evaluation, the integrated system

TABLE VI
LATENCY OF DIFFERENT COMPONENTS OF THE PROPOSED TRACKING FRAMEWORK IN MILLISECONDS ON A NVIDIA A100 GPU. * SHOWS A VERSION WITHOUT TRACK-EMBEDDINGS, \clubsuit USES STREAMPETR [40] WITH A LMM INSTEAD OF DETR3D AS DETECTION TRANSFORMER.

Name	Total	Backbone	Transformer	LMM
STAR-TRACK*	145.6	33.7	80.1	2.2
STAR-TRACK	147.8	33.7	80.1	4.4
StreamPETR \clubsuit	49.7	19.1	13.6	3.3

demonstrated significant improvements in all relevant tracking metrics. Increased track consistency was observed as a particular strength evident from significantly decreased identity switches and track fragmentations.

We hope that this work serves as a foundation for future tracking-by-attention research with the aim of integrating model-based assumptions to end-to-end tracking approaches. While the potential of this has been clearly demonstrated in this work, limitations and opportunities for improvement have also been identified.

Limitations: The implicit association used in the tracking-by-attention scheme falls short in cases with poor motion estimates, since the resulting prediction might impair the re-identification performance in the next frame. This could lead to errors in object position or track losses. In future work, multi-hypothesis tracking [41] could be adopted to model uncertainty in object dynamics and to relax the one-to-one relation of track queries between frames. Additionally, the implicit assignment results in a discrepancy between training and inference time, since the ground truth matching only assigns the correct ground truth object to a single query during training. This could be solved with a non-strict matching approach as demonstrated in 2D tracking [42]. The novel idea of track embeddings is a promising research direction that could be extended to model the uncertainty distribution of each tracked object explicitly.

REFERENCES

- [1] Y. Li, Y. Chen, X. Qi, Z. Li, J. Sun, and J. Jia, "Unifying Voxel-based Representation with Transformer for 3D Object Detection," in *NIPS*, 2022.
- [2] C. Badue, R. Guidolini, R. V. Carneiro, P. Azevedo, V. B. Cardoso, A. Forechi, L. Jesus, R. Berriel, T. M. Paixao, F. Mutz, *et al.*, "Self-Driving Cars: A Survey," *Expert Systems with Applications*, 2021.
- [3] A. Ess, K. Schindler, B. Leibe, and L. Van Gool, "Object Detection and Tracking for Autonomous Navigation in Dynamic Environments," *IJRR*, 2010.
- [4] Y. Wang, V. C. Guizilini, T. Zhang, Y. Wang, H. Zhao, and J. Solomon, "DETR3D: 3D Object Detection from Multi-view Images via 3D-to-2D Queries," in *CoRL*, 2022.
- [5] S. Doll, R. Schulz, L. Schneider, V. Benzin, M. Enzweiler, and H. P. Lensch, "SpatialDETR: Robust Scalable Transformer-Based 3D Object Detection from Multi-View Camera Images with Global Cross-Sensor Attention," in *ECCV*, 2022.
- [6] S. Wang, X. Jiang, and Y. Li, "Focal-PETR: Embracing Foreground for Efficient Multi-Camera 3D Object Detection," *arXiv.org*, vol. arXiv:2212.05505, 2022.
- [7] "nuScenes Tracking Task," <https://nusenes.org/tracking>, accessed: 22.02.23.
- [8] P. Karkus, B. Ivanovic, S. Mannor, and M. Pavone, "DiffStack: A Differentiable and Modular Control Stack for Autonomous Vehicles," in *CoRL*, 2022.
- [9] J. Gu, C. Hu, T. Zhang, X. Chen, Y. Wang, Y. Wang, and H. Zhao, "ViP3D: End-to-end Visual Trajectory Prediction via 3D Agent Queries," *CVPR*, 2023.
- [10] T. Meinhardt, A. Kirillov, L. Leal-Taixe, and C. Feichtenhofer, "TrackFormer: Multi-Object Tracking with Transformers," in *CVPR*, 2022.
- [11] F. Zeng, B. Dong, Y. Zhang, T. Wang, X. Zhang, and Y. Wei, "MOTR: End-to-End Multiple-Object Tracking with Transformer," in *ECCV*, 2022.
- [12] T. Zhang, X. Chen, Y. Wang, Y. Wang, and H. Zhao, "MUTR3D: A Multi-camera Tracking Framework via 3D-to-2D Queries," in *CVPR Workshops*, 2022.
- [13] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention Is All You Need," in *NIPS*, 2017.
- [14] Z. Pang, Z. Li, and N. Wang, "SimpleTrack: Understanding and Rethinking 3D Multi-object Tracking," in *ECCV*, 2023.
- [15] F. Ruppel, F. Faion, C. Gläser, and K. Dietmayer, "Transformers for Multi-Object Tracking on Point Clouds," in *IV*, 2022.
- [16] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liong, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom, "nuScenes: A multimodal dataset for autonomous driving," in *CVPR*, 2020.
- [17] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-End Object Detection with Transformers," in *ECCV*, 2020.
- [18] X. Bai, Z. Hu, X. Zhu, Q. Huang, Y. Chen, H. Fu, and C.-L. Tai, "TransFusion: Robust LiDAR-Camera Fusion for 3D Object Detection with Transformers," in *CVPR*, 2022.
- [19] G. K. Erabati and H. Araujo, "Li3DeTr: A LiDAR based 3D Detection Transformer," in *WACV*, 2023.
- [20] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft COCO: Common Objects in Context," in *ECCV*, 2014.
- [21] T. Yin, X. Zhou, and P. Krahenbuhl, "Center-based 3D Object Detection and Tracking," in *CVPR*, 2021.
- [22] J. Pang, L. Qiu, X. Li, H. Chen, Q. Li, T. Darrell, and F. Yu, "Quasi-Dense Similarity Learning for Multiple Object Tracking," in *CVPR*, 2021.
- [23] E. Ristani and C. Tomasi, "Features for Multi-Target Multi-Camera Tracking and Re-Identification," in *CVPR*, 2018.
- [24] R. Schubert, E. Richter, and G. Wanielik, "Comparison and evaluation of advanced motion models for vehicle tracking," in *IEEE International Conference on Information Fusion*, 2008.
- [25] P. Sun, J. Cao, Y. Jiang, R. Zhang, E. Xie, Z. Yuan, C. Wang, and P. Luo, "TransTrack: Multiple Object Tracking with Transformer," *arXiv.org*, vol. arXiv:2012.15460, 2020.
- [26] D. Ha, A. Dai, and Q. V. Le, "HyperNetworks," *ICLR*, 2017.
- [27] R. E. Kalman, "A New Approach to Linear Filtering and Prediction Problems," *Transactions of the ASME—Journal of Basic Engineering*, 1960.
- [28] W. Luo, J. Xing, A. Milan, X. Zhang, W. Liu, and T.-K. Kim, "Multiple Object Tracking: A Literature Review," *AI*, 2021.
- [29] G. Ciaparrone, F. L. Sánchez, S. Tabik, L. Troiano, R. Tagliaferri, and F. Herrera, "Deep learning in video multi-object tracking: A survey," *Neurocomputing*, 2020.
- [30] Y. Zhou, C. Barnes, J. Lu, J. Yang, and H. Li, "On the Continuity of Rotation Representations in Neural Networks," in *CVPR*, 2019.
- [31] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal Loss for Dense Object Detection," in *CVPR*, 2017.
- [32] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," in *CVPR*, 2016.
- [33] I. Loshchilov and F. Hutter, "Decoupled Weight Decay Regularization," *ICLR*, 2019.
- [34] Z. Pang, Z. Li, and N. Wang, "Simpletrack: Understanding and rethinking 3d multi-object tracking," in *European Conference on Computer Vision*. Springer, 2022, pp. 680–696.
- [35] "MUTR3D Evaluation Github issue #15," <https://github.com/a1600012888/MUTR3D/issues/15>.
- [36] T. Fischer, Y.-H. Yang, S. Kumar, M. Sun, and F. Yu, "Cc-3dt: Panoramic 3d object tracking via cross-camera fusion," *CoRL*, vol. arXiv:2212.01247, 2022.
- [37] Z. Pang, J. Li, P. Tokmakov, D. Chen, S. Zagoruyko, and Y.-X. Wang, "Standing Between Past and Future: Spatio-Temporal Modeling for Multi-Camera 3D Multi-Object Tracking," *arXiv.org*, vol. arXiv:2302.03802, 2023.
- [38] Y. Liu, T. Wang, X. Zhang, and J. Sun, "PETR: Position Embedding Transformation for Multi-View 3D Object Detection," in *ECCV*, 2022.
- [39] Y. Lee, J.-w. Hwang, S. Lee, Y. Bae, and J. Park, "An Energy and GPU-Computation Efficient Backbone Network for Real-Time Object Detection," in *CVPR Workshops*, 2019.
- [40] S. Wang, Y. Liu, T. Wang, Y. Li, and X. Zhang, "Exploring object-centric temporal modeling for efficient multi-view 3d object detection," *arXiv.org*, vol. arXiv:2303.11926, 2023.
- [41] C. Kim, F. Li, A. Ciptadi, and J. M. Rehg, "Multiple Hypothesis Tracking Revisited," in *ICCV*, 2015.
- [42] D. Jia, Y. Yuan, H. He, X. Wu, H. Yu, W. Lin, L. Sun, C. Zhang, and H. Hu, "DETRs with Hybrid Matching," *CVPR*, 2023.