

Convolutional Unscented Kalman Filter for Multi-Object Tracking with Outliers

Shiqi Liu, Wenhan Cao, Chang Liu, Tianyi Zhang, Shengbo Eben Li

Abstract—Multi-object tracking (MOT) is an essential technique for navigation in autonomous driving. In tracking-by-detection systems, biases, false positives, and misses, which are referred to as outliers, are inevitable due to complex traffic scenarios. Recent tracking methods are based on filtering algorithms that overlook these outliers, leading to reduced tracking accuracy or even loss of the object’s trajectory. To handle this challenge, we adopt a probabilistic perspective, regarding the generation of outliers as misspecification between the actual distribution of measurement data and the nominal measurement model used for filtering. We further demonstrate that, by designing a convolutional operation, we can mitigate this misspecification. Incorporating this operation into the widely used unscented Kalman filter (UKF) in commonly adopted tracking algorithms, we derive a variant of the UKF that is robust to outliers, called the convolutional UKF (ConvUKF). We show that ConvUKF maintains the Gaussian conjugate property, thus allowing for real-time tracking. We also prove that ConvUKF has a bounded tracking error in the presence of outliers, which implies robust stability. The experimental results on the KITTI and nuScenes datasets show improved accuracy compared to representative baseline algorithms for MOT tasks.

Index Terms—Multi-object tracking, unscented Kalman filter, outliers

I. INTRODUCTION

MULTI-OBJECT tracking (MOT) is an essential technology for autonomous driving [1]–[3]. This technique provides dynamic environment information for decision-making by continuously monitoring the motion of surrounding objects. Currently, the majority of MOT algorithms adhere to the tracking-by-detection (TBD) paradigm [4]–[8], which consists of three phases: object detection, data association, and filtering. The object detection phase uses a deep neural network to process input data from sensors, such as LiDAR and cameras, to identify objects and determine their motion states, represented as bounding boxes. Subsequently, the data association phase matches these detected bounding boxes with the predicted objects’ trajectories, ensuring the continuous

tracking of the same objects. Lastly, the filtering phase uses the matched detected bounding boxes as measurements to estimate the object’s motion.

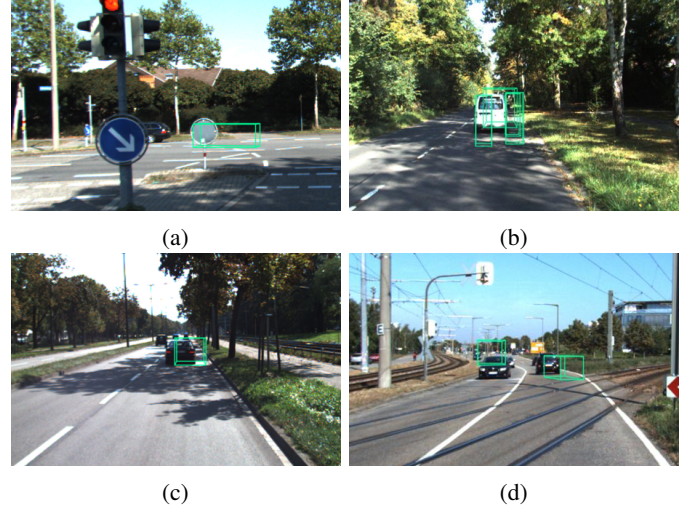


Fig. 1: Examples of detection outliers from LiDAR dataset of KITTI [9]. For better visualization, we project the detection results onto the image instead of using the original point cloud. Note that the bounding boxes are generated by running the PointRCNN detection algorithm [10]. (a) Detection misses: a black car is obscured by a signboard. (b) False positives: redundant bounding boxes on the white car. (c)(d) Detection bias: the bounding boxes deviate from the desired car position.

The accuracy of MOT is significantly influenced by the matched bounding boxes from object detection and data association [11], [12]. Practically, the complexity of cluttered environments and the limitation of detection and association algorithms often lead to a considerable number of biases, false positives, and false negatives (also called misses) [13], [14]. These occurrences lead to outliers, defined as the matched bounding boxes that differ significantly from others for the same objects [15]. An illustration of outliers in real-world scenarios is provided in Fig 1. As input to the filtering phase, these outliers are challenging to model, inevitably causing misspecification of the nominal model and the actual measurement data, which can further deteriorate filtering performance.

Nonetheless, recent tracking methods are rooted in filtering algorithms that overlook the effects of outliers. For example, the most canonical and popular MOT algorithm, known as ABMOT3D [4], is based on the standard Kalman filter (KF). The standard KF naturally assumes that the underlying systems

All correspondence should be sent to Shengbo Eben Li.

Shiqi Liu is with the School of Vehicle and Mobility, Tsinghua University, Beijing, 100084, China. Email: lsq23@mails.tsinghua.edu.cn.

Wenhan Cao is with the State Key Laboratory of Intelligent Green Vehicle and Mobility, Tsinghua University, Beijing, 100084, China. Email: cwh19@mails.tsinghua.edu.cn.

Chang Liu is with the Department of Advanced Manufacturing and Robotics, Peking University, Beijing 100871, China. Email: changliu-coe@pku.edu.cn.

Tianyi Zhang is with the School of Vehicle and Mobility, Tsinghua University, Beijing, 100084, China. Email: 19241041@buaa.edu.cn.

Shengbo Eben Li is with the School of Vehicle and Mobility and College of Artificial Intelligence, Tsinghua University, Beijing, 100084, China. Email: lisb04@gmail.com.

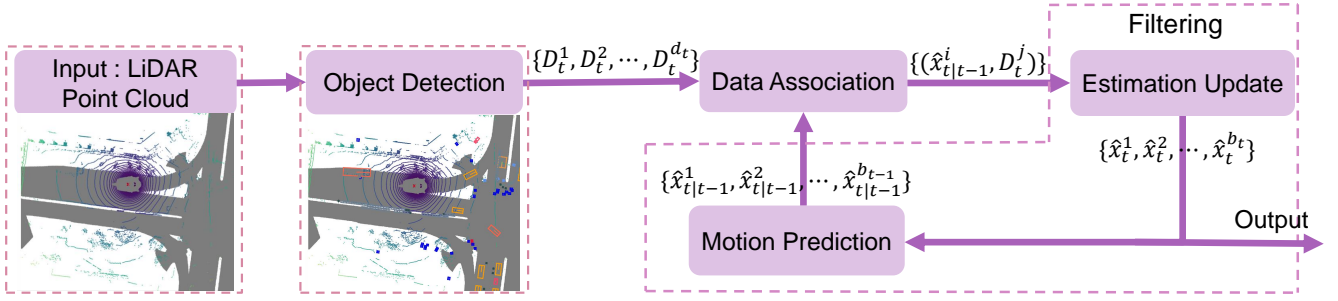


Fig. 2: TBD framework. Firstly, the object detection module detects bounding boxes $\{D_t^1, D_t^2, \dots, D_t^{d_t}\}$ from the raw LiDAR point cloud, providing initial information about the objects in the environment. Next, the motion prediction in filtering module employs a motion model to predict the motions of detected objects, generating $\{\hat{x}_{t|t-1}^1, \hat{x}_{t|t-1}^2, \dots, \hat{x}_{t|t-1}^{b_{t-1}}\}$. Subsequently, the data association module matches these predictions with the current detections for the continuous tracking of the same objects $\{(\hat{x}_{t|t-1}^i, D_t^j)\}$. Finally, the estimation update in filtering module computes the motion estimation of objects $\{\hat{x}_t^1, \hat{x}_t^2, \dots, \hat{x}_t^{b_t}\}$ at timestamp t . As the next frame, object detection will process the incoming sensor data again.

are linear and Gaussian, implying that there are no outliers in the measurement data. Following this line of research, recent works have extended this algorithm with more accurate non-linear models for filtering to improve tracking performance, using extended KF [16] and unscented KF [17]. Unfortunately, these filtering techniques also assume Gaussian distributions for noise. Under such conditions, outliers can cause deviations from Gaussian assumptions, leading to reduced accuracy and even the loss of tracking in MOT. In an effort to tackle time-varying noises, the application of the adaptive cubature Kalman filter in MOT has been proposed [5]. However, the essence of this method is to adapt parameters to time-varying environments rather than to improve the algorithm's robustness under the contamination of outliers.

In this paper, we show that the matched bounding boxes in the TBD framework inevitably contain outliers, which would deteriorate the tracking performance when standard filtering algorithms are applied. In contrast to previous works that aim to avoid outliers by improving the robustness of detection [18], [19] and data association [7], [20], we focus on improving the filtering algorithm's robustness under the contamination of outliers. We show that the misspecification between the nominal model and the measurement data caused by outliers can be quantitatively described using a threshold random variable to capture the uncertainty gap between them. We further demonstrate that this gap can be effectively incorporated through a specialized integral operation akin to convolution. Specifically, this operation aggregates the uncertainty gap throughout the filtering procedure, resulting in a filtering scheme that assigns lower weights to measurement outliers, thereby enhancing the robustness of the filtering. Our contributions are summarized as follows:

- We demonstrate that the outliers can be modeled by introducing an uncertainty gap in the filtering model and mitigated by performing the so-called convolution operation in filtering algorithm. Applying this operation to the popular UKF algorithm results in a novel algorithm, which we refer to as the convolutional UKF (ConvUKF).
- A stability proof for ConvUKF for nonlinear systems with measurement outliers is presented, showing that

the estimation error is bounded in the mean square. Specifically, the upper bound of estimation error is proved to exhibit a linear positive correlation with both the initial state error and the covariance of noise and outliers.

- The effectiveness of the proposed ConvUKF is validated using the real-world KITTI and nuScenes [21] datasets, which demonstrates an improvement in accuracy compared with baseline algorithms. Furthermore, our proposed method has the same computational complexity as UKF because it preserves the Gaussian conjugate structure, which enables real-time tracking.

The structure of the remaining sections is structured as follows: Section II presents the preliminaries of TBD framework and problem formulation of filtering. In Section III, the ConvUKF algorithm and its stability analysis are proposed. In Section IV, we provide an evaluation on the KITTI and nuScenes datasets, along with comparisons with existing filtering methods. Finally, conclusions and limitation discussion are drawn in Section V.

Notation: The notation $M \geq N$ indicates that the matrix $M - N$ is positive semi-definite, while $M > N$ indicates that the matrix $M - N$ is positive definite. For vectors, $\|x\|$ represents the l_2 -norm of the vector x . For matrices, $\|A\|$ means the Frobenius norm of the matrix A . Besides, $I_{m \times n}$ and $0_{m \times n}$ are the identity matrix and zero matrix with dimension $m \times n$.

II. PRELIMINARIES AND PROBLEM FORMULATION

The objective of MOT is to track surrounding objects using raw sensor data. The most applied TBD framework, illustrated in Fig 2, comprises three modules: object detection, data association, and filtering [4]–[6], [8]. In this section, we introduce the preliminaries of the TBD framework and present the problem formulation of filtering under this framework.

A. Preliminaries

Object detection. In MOT, object detection is aimed at identifying and locating the vehicle's surroundings. The process begins with raw data input from onboard sensors like

cameras and LiDAR, which capture the scene around the vehicle. The raw data is then fed into a neural network that extracts features and identifies objects [14]. The output of this module is a set of detected bounding boxes represented as $D_t = \{D_t^1, D_t^2, \dots, D_t^{d_t}\}$, where d_t denotes the number of detections. Each detection, D_t^i is described as a tuple $(p_x, p_y, p_z, \phi, l, w, h)$ for $i \in \{1, 2, \dots, d_t\}$, which includes the position of geometric center (p_x, p_y, p_z) , size of detected bounding boxes (l, w, h) , and yaw angle (ϕ) .

Data association. The data association module aligns the predicted trajectories $\{\hat{x}_{t|t-1}^1, \hat{x}_{t|t-1}^2, \dots, \hat{x}_{t|t-1}^{b_{t-1}}\}$ provided by motion prediction part in filtering module, with the bounding boxes D_t , sourced from the detection module to get the appropriate match $\{(\hat{x}_{t|t-1}^i, D_t^j)\}$, $i \in [1, b_{t-1}]$, $j \in [1, d_t]$, as illustrated by Fig 2. Here b_{t-1} is the number of predicted trajectories in the filtering module. To accomplish the best match, this module encompasses the similarity a_{ij} between each pair of predicted trajectory $\hat{x}_{t|t-1}^i$ and detection D_t^j :

$$a_{ij} = \frac{\|\hat{x}_{t|t-1}^i \cap D_t^j\|}{\|\hat{x}_{t|t-1}^i \cup D_t^j\|}, \forall i \in [1, b_{t-1}], j \in [1, d_t],$$

where $\|\hat{x}_{t|t-1}^i \cap D_t^j\|$ and $\|\hat{x}_{t|t-1}^i \cup D_t^j\|$ represent the intersection volume and the sum volume between the detected objects and the predicted objects, respectively. Using this kind of similarity, the graph-matching problem be formulated and resolved in polynomial time using the Hungarian algorithm [22].

Filtering. The filtering module employs filtering algorithms to estimate objects' motion. Given the necessity of computational efficiency for real-time tracking, the majority of methods employ the Kalman filter family [4], [5], [17], [23], which consists of two parts: motion prediction and estimation update. The motion prediction part utilizes the estimated state $\{\hat{x}_{t-1}^1, \hat{x}_{t-1}^2, \dots, \hat{x}_{t-1}^{b_{t-1}}\}$ from the previous frame ($t-1$) to forecast the movements of objects, resulting in predictions denoted as $\{\hat{x}_{t|t-1}^1, \hat{x}_{t|t-1}^2, \dots, \hat{x}_{t|t-1}^{b_{t-1}}\}$. These predictions will be fed into the data association module for further processing. On the other hand, the estimation update part receives matched pairings $\{(\hat{x}_{t|t-1}^i, D_t^j)\}$ from the data association module. With this input, it calculates the final motion estimations of objects $\{\hat{x}_t^1, \hat{x}_t^2, \dots, \hat{x}_t^{b_t}\}$ through the filtering algorithms.

B. Problem Formulation of Filtering

In MOT, the objective of filtering is to mitigate noise and yield an accurate estimation of the object's motion state $x = [p_x, p_y, p_z, \phi, l, w, h, v_h, v_v, \dot{v}_h, \dot{\phi}]^\top$, including position (p_x, p_y, p_z) , size (l, w, h) , horizontal velocity and vertical velocity (v_h, v_v) , horizontal acceleration \dot{v}_h , yaw angle ϕ , yaw rate $\dot{\phi}$. The horizontal velocity v_h and the horizontal acceleration \dot{v}_h are defined in the direction of the object's orientation, while the vertical velocity v_v is along the z-axis. The system is constructed as a state-space model:

$$x_{t+1} = f(x_t) + \xi_t, \quad (1a)$$

$$y_t = h(x_t) + \zeta_t, \quad (1b)$$

where $f(\cdot)$ and $h(\cdot)$ represent the state transition and measurement models, respectively. The term $x_t \in \mathcal{X} \subseteq R^n$ is defined as the object's motion state while $y_t \in \mathcal{Y} \subseteq R^m$ denotes the specific matched bounding boxes D_t . Furthermore, $\xi_t \sim \mathcal{N}(0, Q_t)$ represents the Gaussian transition noise, characterized by a covariance matrix Q_t . The measurement noise ζ_t , under ideal circumstances, is assumed to follow a Gaussian distribution with $\zeta_t \sim \mathcal{N}(0, R_t)$ with the covariance R_t .

The transition model $f(\cdot)$ is to predict the motion of objects given the state x_t . For motion tracking, there are typically choices, namely the constant velocity (CV) model, the constant turn rate and velocity (CTRV) model, and the constant turn rate and acceleration (CTRA) model [24]. Among them, the CTRA model offers a most precise depiction of motion in real-world object-tracking scenarios by incorporating both turning and acceleration factors [24], whose formulation is shown as

$$\begin{aligned} f(x_t) &= x_t \\ &+ [\Delta p_{x,t}, \Delta p_{y,t}, v_{v,t}\Delta t, \dot{\phi}_t\Delta t, 0_{1 \times 3}, \dot{v}_{h,t}\Delta t, 0_{1 \times 3}]^\top, \\ \Delta p_{x,t} &= \frac{1}{\dot{\phi}_t^2} [(v_{h,t} + \dot{v}_{h,t}\Delta t)\dot{\phi}_t \sin(\phi_t + \dot{\phi}_t\Delta t) \\ &- v_{h,t}\dot{\phi}_t \sin \phi_t + \dot{v}_{h,t} \cos(\phi_t + \dot{\phi}_t\Delta t) - \dot{v}_{h,t} \cos \phi_t], \\ \Delta p_{y,t} &= \frac{1}{\dot{\phi}_t^2} [(-v_{h,t} - \dot{v}_{h,t}\Delta t)\dot{\phi}_t \cos(\phi_t + \dot{\phi}_t\Delta t) \\ &+ v_{h,t}\dot{\phi}_t \cos \phi_t + \dot{v}_{h,t} \sin(\phi_t + \dot{\phi}_t\Delta t) - \dot{v}_{h,t} \sin \phi_t], \end{aligned}$$

where Δt is the sample time. As we introduced in Section II-A, measurements are derived from the matched bounding boxes, represented as $(p_x, p_y, p_z, \phi, l, w, h)$, thus the related measurement model can be formalized as

$$h(x_t) = [I_{7 \times 7}, 0_{4 \times 4}] x_t.$$

Remark 1. The state space model (SSM) (1) can also be represented as the hidden Markov model (HMM):

$$x_t \sim p(x_t|x_{t-1}), y_t \sim p(y_t|x_t), \quad (2)$$

where $p(x_t|x_{t-1})$ represents the transition probability, and $p(y_t|x_t)$ is the likelihood probability. Under the ideal conditions, we have $p(x_t|x_{t-1}) = \mathcal{N}(x_t; f(x_{t-1}), Q_t)$ and $p(y_t|x_t) = \mathcal{N}(y_t; h(x_t), R_t)$.

III. CONVOLUTIONAL UNSCENTED KALMAN FILTER

The UKF, as one of the Kalman family filters presented by [25], stands out as a widely adopted method for its efficiency and high precision in estimating the state of nonlinear motion in MOT [17]. However, the UKF assumes an ideal HMM as described in (2), and the presence of outliers in measurements can significantly challenge its estimation performance. In this section, we analyze the HMM with outliers and present the ConvUKF algorithm under this analysis. Finally, we demonstrate the stability of ConvUKF in nonlinear systems with outliers.

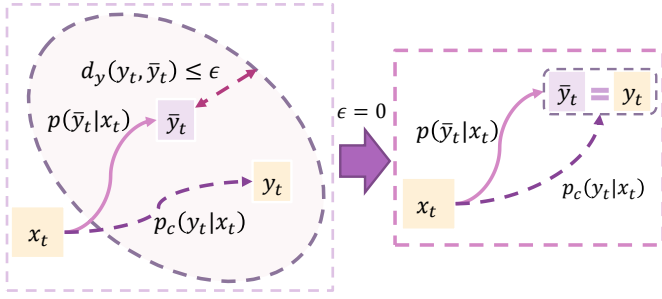


Fig. 3: Illustration of the HMM with outliers and without outliers. The nominal likelihood probability in HMM projects the state in the previous time, denoted as x_{t-1} , to a modeled virtual measurement \bar{y}_t . However, due to stochastic nature of outliers, there is an uncertainty gap between y_t and \bar{y}_t . We assume their gap distance $d_y(y_t, \bar{y}_t)$ can be modeled by a stochastic inequality ϵ , i.e., $d_y(y_t, \bar{y}_t) \leq \epsilon$. Especially when $\epsilon = 0$, meaning no outliers, the system is reduced to the ideal scene $\bar{y}_t = y_t$.

A. Algorithm Design

In MOT, attaining precise model $p(y_t|x_t)$ in (2) is often unfeasible due to the inevitable biases, false positives, and false negatives [13], [14]. Consequently, it is necessary to differentiate the actual measurement variable y_t , often contaminated by outliers, and the virtual measurement variable \bar{y}_t , which is shown in Fig 3. The former is an accurate yet unattainable description of the measurement, while the latter is an artificial construct generated by nominal models.

Following research [26], we construct a stochastic inequality to model the uncertainty gap between y_t and \bar{y}_t by a threshold random variable ϵ , represented as

$$d_y(y_t, \bar{y}_t) \leq \epsilon \quad (3)$$

Here, $d_y : \mathcal{Y} \times \mathcal{Y} \rightarrow R$ denotes the distance function. Note that the cumulative distribution function of the threshold random variable is assumed to be known and denoted as F_ϵ . Combined with the original HMM (2), this constitutes the HMM with outliers:

$$x_t \sim p(x_t|x_{t-1}), \bar{y}_t \sim p(\bar{y}_t|x_t), d_y(y_t, \bar{y}_t) \leq \epsilon, \quad (4)$$

where $p(\bar{y}_t|x_t) = \mathcal{N}(\bar{y}_t; h(x_t), R_t)$ is the likelihood probability for the virtual measurement variable \bar{y}_t under the ideal condition, called the nominal likelihood. Note that when $\epsilon = 0$, (4) can be reduced to the ideal HMM without outliers, as illustrated by Fig 3. Due to the fact that the physical world and the modeling of a system are mutually exclusive at any given moment, it is reasonable to assume that y_t and \bar{y}_t are conditionally independent given x_t , and ϵ is also independent of both y_t and \bar{y}_t [26].

Under the formulation of HMM with outliers in (4), we can calculate the so-called convolutional likelihood probability [26] by additionally conditioning on stochastic inequality (3):

$$p(y_t|x_t, d(y_t, \bar{y}_t) \leq \epsilon) \propto \int_{\bar{y}_t} (1 - F_\epsilon(d_y(y_t, \bar{y}_t))) p(\bar{y}_t|x_t) d\bar{y}_t. \quad (5)$$

Here, $p_c(y_t|x_t) := p(y_t|x_t, d(y_t, \bar{y}_t) \leq \epsilon)$ is the convolutional likelihood. The convolutional likelihood owes its name as it is calculated by an integral operation akin to convolution, which integrates the uncertainty gap ϵ caused by outliers. Specifically, $1 - F_\epsilon(d_y(y_t, \bar{y}_t))$ serves as the kernel function, which is a weighting coefficient of the nominal likelihood $p(\bar{y}_t|x_t)$ based on the distance $d_y(y_t, \bar{y}_t)$. The convolutional likelihood is proved to be a more informative estimate of the actual likelihood function than the nominal one [26] when the system is contaminated by outliers.

A practical challenge with this scheme is whether this newly defined likelihood has an analytical form, given that the convolution operation typically lacks a closed-form solution [27]. Fortunately, for the common quadratic-form distance, if the threshold random variable follows an exponential distribution, the nominal model with Gaussian noise possesses an analytical form of the convolutional likelihood function. This is demonstrated in the subsequent lemma.

Lemma 1 ([26]). *Consider the following nominal system model $p(\bar{y}_t|x_t) = \mathcal{N}(\bar{y}_t; h(x_t), R_t)$. If $d_y(y, \bar{y}) = \|y - \bar{y}\|^2$ and $\epsilon \sim \text{Exp}(\gamma)$ with $\gamma > 0$ as the parameter of the exponential distribution, then we have $p_c(y_t|x_t) = \mathcal{N}(y_t; g(x_t), R_t + \frac{1}{2\gamma} \cdot I_{m \times m})$.*

As shown in this lemma, if the nominal likelihood function is Gaussian, its convolutional counterpart is still Gaussian. This property allows us to preserve the original Gaussian conjugate structure if we apply the convolutional update to the well-known KF filter family to which UKF belongs. By reshaping the original likelihood of UKF to a convolutional likelihood, we have the following ConvUKF algorithm:

- **Initialization.** Given: $\hat{x}_0, \hat{P}_0, Q, R, \gamma, a$, repeat steps 1–3 for $t = 0, 1, 2, \dots$
- **Step 1: Select sigma points.** Here are the Julier sigma points: $\mathcal{X}_{0,t} = \hat{x}_t; \mathcal{X}_{i,t} = \hat{x}_t + (a\sqrt{n\hat{P}_t})_i; \mathcal{X}_{i,t} = \hat{x}_t - (a\sqrt{n\hat{P}_t})_i, i = 1, 2, \dots, n$, where a is a proportion parameter and $(a\sqrt{n\hat{P}_t})_i$ is the vector of the i th column of the matrix square root.
- **Step 2: Prediction.**

$$\begin{aligned} \hat{x}_{t+1|t} &= \sum_{i=0}^{2n} \omega_i f(\mathcal{X}_{i,t}), i = 0, 1, 2, \dots, 2n, \\ \check{\mathcal{X}}_{i,t+1|t} &= f(\mathcal{X}_{i,t}) - \hat{x}_{t+1|t}, \\ \hat{P}_{t+1|t} &= \sum_{i=0}^{2n} \omega_i \check{\mathcal{X}}_{i,t+1|t} \check{\mathcal{X}}_{i,t+1|t}^\top + Q_t, \end{aligned} \quad (6)$$

where $\omega_0 = 1 - \frac{1}{a^2}, \omega_i = \frac{1}{2na^2}, i = 1, 2, \dots, 2n$ are the scalar weights with $\sum_{i=0}^{2n} \omega_i = 1$. It is noticeable that $\hat{x}_{t+1|t}$ is the prediction state, and $\hat{P}_{t+1|t}$ is the predicted covariance at time step t . Both $\hat{x}_{t+1|t}$ and $\hat{P}_{t+1|t}$ will be corrected in the update step.

- **Step 3: Update.**

$$\begin{aligned} \hat{y}_{t+1} &= \sum_{i=0}^{2n} \omega_i h(\mathcal{X}_{i,t}), i = 0, 1, 2, \dots, 2n, \\ \check{\mathcal{Y}}_{i,t+1} &= h(\mathcal{X}_{i,t}) - \hat{y}_{t+1}, \end{aligned}$$

$$\hat{P}_{yy,t+1} = \sum_{i=0}^{2n} \omega_i \check{\mathcal{Y}}_{i,t+1} \check{\mathcal{Y}}_{i,t+1}^\top + R_{t+1} + \frac{1}{2\gamma} \cdot I_{m \times m}, \quad (7a)$$

$$\hat{P}_{xy,t+1} = \sum_{i=0}^{2n} \omega_i \check{\mathcal{X}}_{i,t+1|t} \check{\mathcal{Y}}_{i,t+1}^\top, \quad (7b)$$

$$K_{t+1} = \hat{P}_{xy,t+1} \hat{P}_{yy,t+1}^{-1}, \quad (7c)$$

$$\hat{x}_{t+1} = \hat{x}_{t+1|t} + K_{t+1}(y_{t+1} - \hat{y}_{t+1}), \quad (7d)$$

$$\hat{P}_{t+1} = \hat{P}_{t+1|t} - K_{t+1} \hat{P}_{xy,t+1}^\top. \quad (7e)$$

In (7a) and (7b) $\hat{P}_{yy,t+1}$ and $\hat{P}_{xy,t+1}$ represent the measurement covariance and the cross-covariance at time step $t+1$. The term $\frac{1}{2\gamma} \cdot I_{m \times m}$ in (7a) is derived from the convolutional probability in Lemma 1. Meanwhile, \hat{x}_{t+1} and \hat{P}_{t+1} denote the estimated state and its associated covariance, respectively.

In MOT, the parameter γ in (7a) is associated with the extent of outlier contamination in measurements, which is challenging to determine.

Inspired by the adaptive filtering algorithms in [5], we adopt a similar adaptive update rule for γ :

$$\gamma^{(t+1)} = (1 - \tau) \gamma^{(t)} + \tau \frac{\gamma^{(t)}}{1 + \exp\{-2\gamma^{(t)}[\exp(-\gamma^{(t)}) - \frac{\tilde{y}_t}{m}]\}}, \quad (8)$$

where τ , set to the empirical value of 0.05 according to [5], is the temperature parameter regulating the update rate, and $\tilde{y}_t = y_{t+1} - \hat{y}_{t+1}$ denotes the measurement error. The underlying concept is straightforward: as the measurement error increases, indicating a larger mismatch between the distribution of measurement data and its corresponding model, γ should be reduced to accommodate this greater uncertainty gap. Conversely, if the error diminishes, the opposite adjustment will be made.

B. Stability Analysis

To employ the proposed algorithm in MOT, it is necessary to ensure filtering stability in nonlinear systems with outliers. Mathematically, filtering stability requires that the estimation error, $\tilde{x}_{t+1} = x_{t+1} - \hat{x}_{t+1}$, is bounded in the mean square, i.e., $\mathbb{E}[\|\tilde{x}_{t+1}\|^2] < +\infty$ [28]. One challenge for stability analysis is that ConvUKF uses UT to approximate the nonlinear system with a linear system, which inevitably introduces approximation error. A viable solution to compensate for this approximation error is to introduce auxiliary variables in the linear approximation [28], [29], leading to the following assumption:

Assumption 1. (Rectified Linearization [28], [29]) For the nonlinear system with outliers in (1), the prediction error $\tilde{x}_{t+1|t} = x_{t+1} - \hat{x}_{t+1|t}$ and measurement error $\tilde{y}_{t+1} = y_{t+1} - \hat{y}_{t+1}$ can be linearly updated with the auxiliary variables $\alpha_t = \text{diag}\{\alpha_{1,t}, \dots, \alpha_{n,t}\}$ and $\beta_t = \text{diag}\{\beta_{1,t}, \dots, \beta_{n,t}\}$.

$$\tilde{x}_{t+1|t} = \alpha_t F_t \tilde{x}_{t|t-1} - \alpha_t F_t K_t \tilde{y}_{t|t-1} + \xi_t, \quad (9a)$$

$$\tilde{y}_{t+1} = \beta_{t+1} H_{t+1} \tilde{x}_{t+1|t} + \zeta_{t+1}, \quad (9b)$$

where $F_t = (\frac{\partial f(x)}{\partial x}|_{x=\hat{x}_t})$ and $H_{t+1} = (\frac{\partial h(x)}{\partial x}|_{x=\hat{x}_{t+1|t}})$ are Jacobian matrices.

Assumption 1 essentially assumes the neglect of higher-order terms when using UT can be rectified by the auxiliary variables α_t and β_t . Based on this assumption, the predicted covariance $\hat{P}_{t+1|t}$ and the measurement covariance $\hat{P}_{yy,t+1}$ can be computed in the following proposition.

Proposition 1. (Rectified Covariance Update [29]) Under Assumption 1, the predicted covariance $\hat{P}_{t+1|t}$ and the measurement covariance $\hat{P}_{yy,t+1}$ are derived as:

$$\begin{aligned} \hat{P}_{t+1|t} &= [\alpha_t F_t (I - K_t \beta_t H_t)] \hat{P}_{t|t-1} \\ &\quad \times [\alpha_t F_t (I - K_t \beta_t H_t)]^\top + Q_t + \Delta P_{t+1|t}, \end{aligned} \quad (10a)$$

$$\begin{aligned} \hat{P}_{yy,t+1} &= (\beta_{t+1} H_{t+1}) \hat{P}_{t+1|t} (\beta_{t+1} H_{t+1})^\top \\ &\quad + R_{t+1} + \frac{1}{2\gamma} \cdot I_{m \times m} + \Delta P_{yy,t+1}, \end{aligned} \quad (10b)$$

where ΔP_{t+1} and $\Delta P_{yy,t+1}$ denote the error covariance, capturing the discrepancies between the actual nonlinear dynamics and their linear approximations, defined as (11).

$$\begin{aligned} \Delta P_{t+1|t} &= \mathbb{E}\{[\alpha_t F_t (I - K_t \beta_t H_t) \tilde{x}_{t+1|t}] \\ &\quad \times [\alpha_t F_t (I - K_t \beta_t H_t) \tilde{x}_{t+1|t}]^\top \\ &\quad - [\alpha_t F_t (I - K_t \beta_t H_t)] \hat{P}_{t|t-1} [\alpha_t F_t (I - K_t \beta_t H_t)] \\ &\quad - P_{t+1|t} + \hat{P}_{t+1|t}, \end{aligned} \quad (11a)$$

$$\begin{aligned} \Delta P_{yy,t+1} &= \mathbb{E}[(\beta_{t+1} H_{t+1}) \tilde{x}_{t+1|t} \tilde{x}_{t+1|t}^\top (\beta_{t+1} H_{t+1})^\top \\ &\quad - (\beta_{t+1} H_{t+1}) \hat{P}_{t+1|t} (\beta_{t+1} H_{t+1})^\top \\ &\quad - P_{yy,t+1} + \hat{P}_{yy,t+1} \end{aligned} \quad (11b)$$

The proof of Proposition 1 can be found in Appendix A. This proposition implies that by adding the covariance error term, we can derive the calculation formula for the estimated covariance. The covariance error term can be seen as the rectification covariance of noise. Considering the presence of outliers, we define the nominal measurement covariance R_t and the discrepancy between the true covariance and the nominal one induced by outliers as ΔR_t . Subsequently, we will perform a stability analysis under the assumption that certain variables within the system adhere to boundedness, as outlined below.

Assumption 2. We assume that the following 3 parts of the hypothesis are satisfied for all $t > 0$.

- **Bounded dynamics:** The dynamics of the system are bounded by real constants ($f_u, h_u, \beta_u, \alpha_u$):

$$F_t F_t^\top \leq f_u^2 I, \quad H_t H_t^\top \leq h_u^2 I, \quad \alpha_t \leq \alpha_u I, \quad \beta_t \leq \beta_u I.$$

- **Bounded noise:** There are real constants ($q_l, q_u, r_u, \Delta r_u$) that bound the noise covariance:

$$q_l I \leq Q_t \leq q_u I, \quad R_t \leq r_u I, \quad 0 \leq \Delta R_t \leq \Delta r_u I.$$

- **Bounded rectified covariance.** For the error covariance, this assumption establishes bounds through real constants $\Delta p_l, \Delta p_u, \Delta p_{yy,u}, p_l, p_u$:

$$\begin{aligned} \Delta p_l I &\leq \Delta P_{t+1|t} \leq \Delta p_u I, \quad \Delta P_{yy,t} \leq \Delta p_{yy,u} I, \\ p_l I &\leq \hat{P}_t \leq p_u I. \end{aligned}$$

TABLE I: PERFORMANCE COMPARISON ON KITTI VALIDATION SET.

Dataset	Method	sAMOTA% ↑	AMOTA% ↑	AMOTP% ↑	MT% ↑	ML% ↓	IFN↓	FRAG↓	Time (ms)↓
KITTI	UKF	76.30	31.47	51.83	47.57	5.41	898	321	0.4225
	HuberUKF	84.42	38.12	61.47	69.19	3.24	725	50	0.9496
	AUKF	84.63	38.02	74.80	70.81	2.16	843	270	0.6510
	IEKF	83.30	37.30	58.31	68.11	4.86	828	69	0.1890
	REKF	89.51	42.15	65.23	71.89	3.78	792	47	0.3225
	ICKF	88.34	41.32	45.50	70.81	3.24	695	67	0.7186
	ConvUKF	93.32	45.46	78.09	75.68	3.78	799	17	0.5726

The best performances are marked with **bold** font. The symbol ↑ indicates better performance for larger values.

TABLE II: PERFORMANCE COMPARISON ON nuSCENES VALIDATION SET.

Dataset	Method	sAMOTA% ↑	AMOTA% ↑	AMOTP% ↑	MT% ↑	ML% ↓	IDS↓	FRAG↓	Time (ms)↓
nuScenes	UKF	62.66	21.93	32.07	37.90	28.29	3973	4740	0.3102
	HuberUKF	66.13	23.69	39.05	39.67	30.09	799	1499	0.6024
	AUKF	67.15	24.19	53.86	41.49	27.76	1959	2520	0.3138
	IEKF	56.11	18.02	25.01	33.46	30.90	662	1817	0.1946
	REKF	64.27	22.54	33.03	39.60	29.19	721	1440	0.2906
	ICKF	66.12	23.50	16.19	38.34	30.46	899	1959	0.7238
	ConvUKF	69.12	25.05	49.58	42.86	29.34	703	1171	0.5788

Assumption 2 is a prerequisite for filtering stability [28]–[30], assuming stable dynamics and finite noise conditions. Based on Assumption 2, we can derive that the bound of the predicted covariance $\hat{P}_{t+1|t}$ and the norm of Kalman gain $\|K_{t+1}\|$.

Proposition 2. *The prediction covariance $\hat{P}_{t+1|t}$ and the norm of the Kalman gain matrix $\|K_{t+1}\|$ are bounded as $\hat{p}_l I \leq \hat{P}_{t+1|t} \leq \hat{p}_u I$ and $\|K_{t+1}\| \leq K_u$ with $\hat{p}_l = p_l, \hat{p}_u = p_u \bar{a}^2 f_u^2 + q_u + \Delta p_u, K_u = \sqrt{n}(p_u \alpha_u^2 f_u^2 + q_u + \Delta p_u)$.*

Proposition 2 provides the bound of $\hat{P}_{t+1|t}$ and $\|K_{t+1}\|$ with proof of this proposition detailed in Appendix B. Then we can formulate and prove the following stability theorem for the ConvUKF:

Theorem 1. *For nonlinear stochastic systems with outliers (1), the estimation error of ConvUKF is bounded in the mean square if Assumption 1 and 2 are satisfied. Specifically, the bound has a linear positive correlation with the initial state error and the upper bound of measurement noise and outliers' covariance.*

$$\mathbb{E}\{\|\tilde{x}_{t+1}\|^2\} \leq C_1 \mathbb{E}\{\|\tilde{x}_0\|^2\} + C_2(r_u + \Delta r_u) + C_3, \quad (12)$$

where C_1, C_2, C_3 are real constants uncorrelated with $\tilde{x}_0, r_u, \Delta r_u, \{y_k\}_{k=1:t+1}$ and $\{\hat{x}_k\}_{k=1:t+1}$.

The proof of Theorem 1 and detailed definition of C_1, C_2, C_3 can be found in Appendix C. The basic proof sketch for Theorem 1 is to construct a quadratic function of the error as a Lyapunov function to demonstrate stability. This technique is widely employed in previous works [28], [29]. Notably, we consider the systems with measurement outliers, a factor overlooked in the aforementioned works.

Remark 2. *From Theorem 1, we can deduce that the bound of estimation error is linear positively correlated with the initial state error and the upper bound of measurement noise and outliers' covariance. Specifically, the impact of various system parameters on the error bound of state estimation can be*

inferred through the definition of constants C_1, C_2 , and C_3 in (26). For example, the estimation error of a system with lower f_u will be less affected by outliers.

IV. EXPERIMENTS

A. Settings

We evaluate our proposed method using the validation set of KITTI and nuScenes datasets. These datasets are widely adopted in the field of autonomous driving for their collection of raw data captured through onboard sensors during real-world driving scenarios. Specifically, we utilize raw LiDAR cloud point as our input data. To valid our method, which is employed in the filtering module, we employ the trained model PointRCNN from [10] (in KITTI dataset) and the Megvii from [31] (in nuScenes dataset) for object detection and the Hungarian algorithm [22] described in Section II-A for data association. For a fair comparison, all baseline algorithms utilize the same object detection and data association methodologies.

To evaluate the efficacy of our method, we conduct a comparative analysis with several established algorithms, including the UKF, Huber unscented Kalman filter (HuberUKF) [32], Adaptive UKF (AUKF) [33], iterative extended Kalman filter (IEKF) [34], outlier-robust extended Kalman filter (REKF) [35] and improved cubature kalman filter (ICKF) [36], on the validation sets of the KITTI and nuScenes. HuberUKF is a classical outlier-robust filtering method that incorporates the Huber loss to form a more robust cost function in the UKF. AUKF is a novelly adaptive filtering method that dynamically adjusts the covariance parameters based on estimation errors to handle complex scenes effectively. IEKF utilizes the concept of Newton's method iteration to improve nonlinear approximation effects. REKF and ICKF adopt factor adjustment, robust weighting, and variance calibration, demonstrating commendable robustness in practical applications. To ensure a fair benchmark comparison, the parameters for each baseline method are configured according to the specifications provided in their respective original literature.

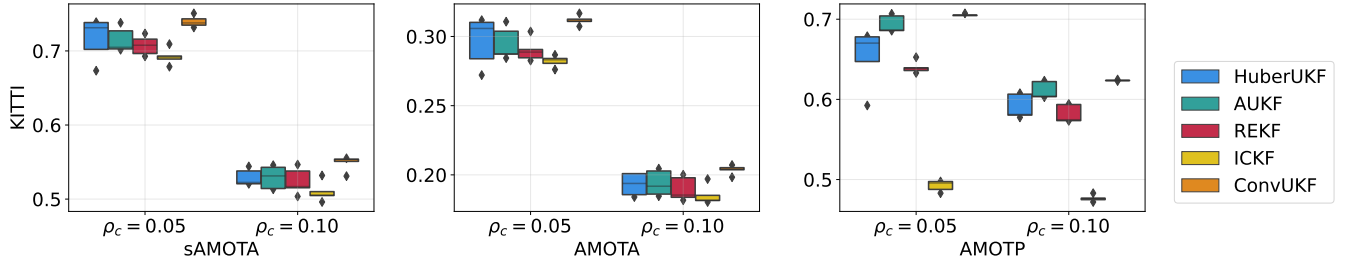


Fig. 4: Boxplot illustrating sAMOTA, AMOTA, and AMOTP values obtained from HuberUKF, AUKF, REKF, ICKF, and ConvUKF with contamination probabilities $\rho_c \in \{5\%, 10\%\}$ in the KITTI validation set. The circles are the data points that fall significantly outside the interquartile range.

TABLE III: THE RESULTS ON THE CONTAMINATED KITTI VALIDATION SET.

Probability	Method	sAMOTA %	AMOTA %	AMOTP %
$\rho_c = 0.05$	HuberUKF	71.66	29.66	65.53
	AUKF	71.69	29.47	69.07
	REKF	70.73	29.00	63.99
	ICKF	69.23	28.21	49.23
	ConvUKF	73.96	31.17	70.53
$\rho_c = 0.10$	HuberUKF	52.92	19.30	59.06
	AUKF	52.95	19.40	61.11
	REKF	52.41	18.98	58.18
	ICKF	50.96	18.51	47.66
	ConvUKF	54.99	20.37	62.37

TABLE IV: THE RESULTS ON THE CONTAMINATED nuSCENES VALIDATION SET.

Probability	Method	sAMOTA %	AMOTA %	AMOTP %
$\rho_c = 0.05$	HuberUKF	61.81	21.01	42.38
	AUKF	55.15	18.20	39.62
	REKF	61.18	20.64	36.33
	ICKF	60.65	20.26	23.04
	ConvUKF	64.37	22.33	48.84
$\rho_c = 0.10$	HuberUKF	58.31	19.05	42.86
	AUKF	54.93	17.72	44.46
	REKF	56.88	18.85	42.26
	ICKF	56.32	17.94	27.54
	ConvUKF	59.97	19.92	47.73

MOT challenge benchmarks employ the following metrics: scaled Average Multi-Object Tracking Accuracy (sAMOTA), Average Multi-Object Tracking Accuracy (AMOTA), Average Multi-Object Tracking Precision (AMOTP), Mostly Tracked (MT), Mostly Lost (ML), Ignored False Negatives (IFN) and Fragmentation (FRAG). For detailed definitions of these metrics, please refer to [4], [37]. As for hardware, our experiments were conducted on a desktop computer equipped with an Intel(R) Core(TM) i9-12900K processor and an NVIDIA 3090 Ti GPU.

B. Results

Results on KITTI Validation Set. The comparative results are detailed in Table I, which shows that ConvUKF algorithms demonstrate superior enhancements across almost all metrics compared to the baseline methods. ConvUKF achieves an impressively high sAMOTA (93.32%) while maintaining the best AMOTA (45.46%), AMOTAP (78.09%), MT (75.68%), and FRAG (17), highlighting the promising potential of integrating ConvUKF into filter-based MOT methods.

Although the ML and IFN of ConvUKF are slightly inferior to those of AUKF and ICKF, ConvUKF obtains better comprehensive performance, representing a 4.26% improvement in sAMOTA compared to the best sAMOTA achieved by REKF.

Results on nuScenes Validation Set. To verify the effectiveness of ConvUKF in different situations, we conducted the same experiment on the nuScenes validation set, as shown in Table II. ConvUKF also demonstrates a better comprehensive performance among all the baseline methods, achieving the best sAMOTA (69.12%), AMOTA (25.05%),

MT (42.86%) and FRAG (1171). AUKF and ICKF achieve the best performance in ML and IDS, respectively, but with a relatively low sAMOTA, reflecting poor overall accuracy. This highlights that while other filtering methods may excel in specific performance indicators, they often fall short in overall tracking accuracy. In contrast, ConvUKF demonstrates a high level of overall performance, achieving strong results across all metrics. We notice that the improvement of ConvUKF's sAMOTA on nuScenes (2.93%) is lower compared with KITTI (10.27%). Our analysis suggests that the nuScenes dataset, with approximately 40,000 frames, contains about three times as much data as the KITTI dataset, which has 15,000 frames, and utilizes a more accurate LiDAR sensor [21]. Consequently, the proportion of outliers in nuScenes is smaller, resulting in a smaller degree of improvement with our method.

Validation for Robustness. To more rigorously validate the robustness of our method, we introduce additional data contamination as outliers to the detection results. This is achieved by masking the bounding boxes at varying probabilities $\rho_c \in \{5\%, 10\%\}$ to simulate scenarios where the object detection fails to recognize targets or detection errors are significant enough to lead to incorrect association matching. Due to the randomness of the masking process, we repeat the experiment 5 times on the validation sets of KITTI and nuScenes. In this experiment, we focus on the three most comprehensive tracking performance indicators: sAMOTA, AMOTA, and AMOTP.

For the KITTI validation set, we present a boxplot of sAMOTA, AMOTA, and AMOTP values in Figure 4 with the average metric values detailed in Table III. HuberUKF,

TABLE V: THE RESULTS OF THE ABLATION EXPERIMENT ON KITTI VALIDATION SET.

Method	sAMOTA%	AMOTA%	AMOTP%
HuberUKF	84.42	38.12	61.47
AUKF	84.63	38.02	74.80
REKF	89.51	42.15	65.23
ICKF	88.34	41.32	45.50
ConvUKF ($\gamma = 1 \times 10^{-1}$)	84.69	38.31	61.35
ConvUKF ($\gamma = 1 \times 10^{-2}$)	91.79	44.38	69.60
ConvUKF ($\gamma = 1 \times 10^{-3}$)	92.70	45.02	72.11
ConvUKF (adaptive)	93.32	45.46	78.09

AUKF, and REKF exhibit similar tracking performance across different levels of outlier contamination. ICKF shows the lowest accuracy but maintains stability with very minimal variance. Notably, ConvUKF consistently demonstrates exceptional stability with a minimal variance while achieving the highest sAMOTA, AMOTA, and AMOTP. In comparison, ConvUKF shows improvements of 3.17% 5.77% 2.11% in sAMOTA, AMOTA, and AMOTP, respectively, at $\rho_c = 0.05$ and 3.86% 5.00% 2.06% at $\rho_c = 0.10$.

For the nuScenes validation set, as shown in Table IV, the performance of AUKF decreases significantly as the contamination probability increases, while HuberUKF, REKF, and ICKF which emphasize robustness, perform better. ConvUKF consistently outperforms other methods in terms of sAMOTA, AMOTA, and AMOTP across all levels of contamination probability, validating its robustness in handling outliers in MOT.

Ablation Studies. We also conduct an ablation experiment on the KITTI dataset with the same settings as described in Section IV-A. The results are presented in Table V. Compared with the conventional ConvUKF using a carefully chosen γ , our ConvUKF with the adaptive update rule demonstrates superior performance across sAMOTA (+0.67%), AMOTA (+0.98%), and AMOTP (+8.28%) metrics. This showcases the impact of leveraging the adaptive parameter trick on the overall performance.

Additionally, most variants of ConvUKF outperform the baseline algorithms, further highlighting the efficacy of ConvUKF in MOT.

Run-time Comparison. We compare the calculation time of filtering in MOT on the KITTI validation set (Table I) and nuScenes validation set (Table II). Since the overall trend and size relationship are consistent across both datasets, we use the KITTI results for specific analysis.

IEKF (0.1890ms) requires the least computation time among all the methods because this method does not need to obtain sigma points, resulting in lesser computation. However, IEKF sacrifices some overall accuracy because the iterative linear approximation does not accurately approximate the non-linear motion of objects. HuberUKF and ICKF, which require both sampling sigma points and Huber Loss calculation, take the longest time to calculate. Due to the Gaussian conjugate structure similar to UKF, ConvUKF (0.5726ms), due to the Gaussian conjugate structure similar to UKF, shows performance comparable to UKF (0.4225ms), making it suitable for real-time object tracking applications.

C. Visualization and Discussion

For an intuitive analysis, we visualize the MOT results on KITTI validation set from 4 continuous frames with an interval of 0.8s. It is noticed that we use the raw LiDAR cloud point as the input but visualize the results on the camera images for better presentation. As the black car in front begins to turn in Fig 5a, the AUKF's tracking (ID:3959) gradually deviates from the car's actual position until it completely loses the target. Conversely, in Fig 5b ConvUKF (ID:3913) consistently maintains accurate tracking of the black car under these conditions.

In both the validation sets of KITTI and nuScenes, as well as in their data-contaminated counterparts, ConvUKF demonstrates strong performance across key metrics such as sAMOTA, AMOTP and so on. Notably, ConvUKF achieves this without significantly extending computational timeframes, ensuring practical feasibility. Moreover, ConvUKF exhibits superior robustness and minimal variance in outlier handling—a pivotal advantage for sustaining the stability and reliability of MOT systems amidst challenging conditions.

V. CONCLUSION AND LIMITATION

In this paper, we formulate a stochastic inequality to characterize the uncertainty gap between the distribution of measurements and its corresponding model for filtering. To mitigate this gap, we introduce a novel convolution-like internal process within the UKF, resulting in the ConvUKF. We demonstrate that ConvUKF not only maintains the conjugate structure of UKF but also obtains more robustness. Furthermore, we prove the stability of ConvUKF in nonlinear systems with outliers, showing that ConvUKF's estimation error is bounded in the mean square sense. Through experiments conducted on real-world datasets (KITTI and nuScenes), we showcase a performance enhancement along with acceptable computation time.

Future work will focus on how to overcome the limitations of our approach. Firstly, identify the accurate parameter γ in ConvUKF instead of relying on the adaptive update rule. We will explore employing ConvUKF in various state-of-the-art MOT algorithms [6]–[8]. Additionally, we plan to apply ConvUKF for MOT in real-world environments, extending beyond our current offline experiments using the dataset.

ACKNOWLEDGMENTS

This study is supported by National Key R&D Program of China with 2022YFB2502901, NSF China under 52221005, Tsinghua University Initiative Scientific Research Program and Tsinghua University-Toyota Joint Research Center for AI Technology of Automated Vehicle.

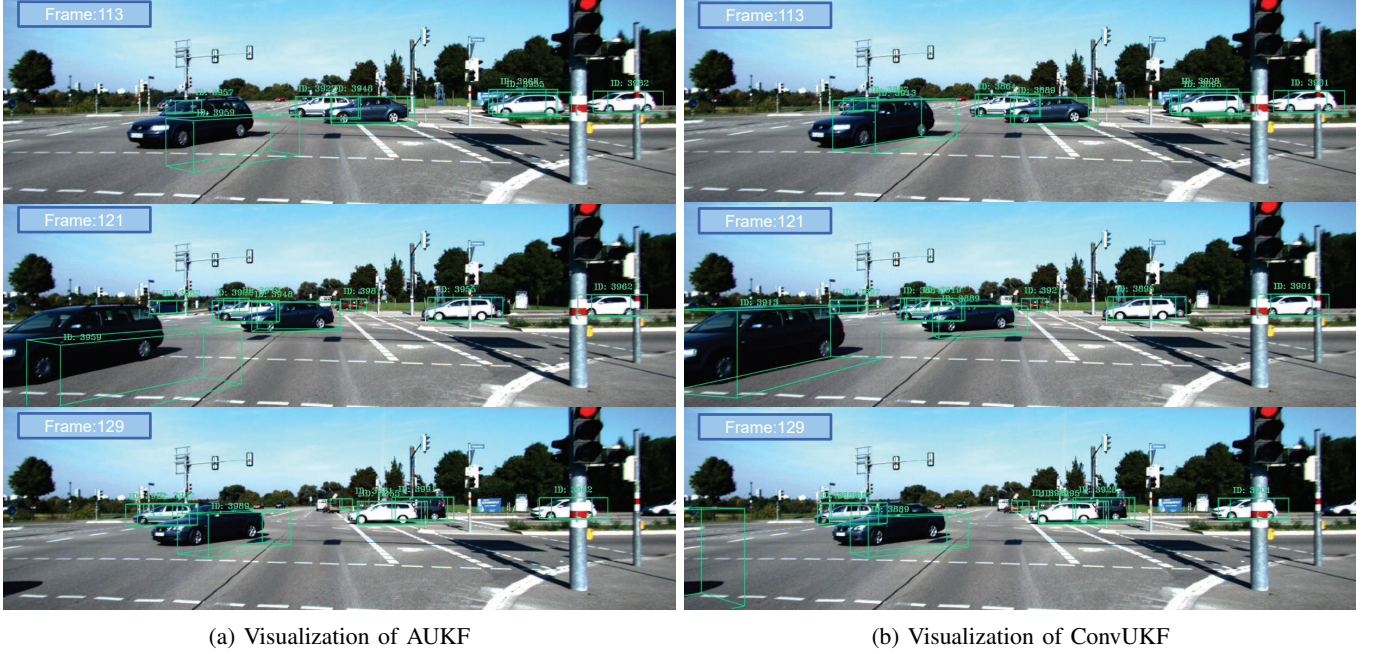


Fig. 5: Visual comparison of MOT results between AUKF (left) and ConvUKF (right), using data from 4 continuous frames with an interval of 0.8s in the KITTI dataset. For better presentation, we visualize the results on the camera images.

APPENDIX A PROOF OF PROPOSITION 1

From (9a), the actual prediction error covariance can be derived as

$$\begin{aligned}
 P_{t+1|t} &= \mathbb{E}[\tilde{x}_{t+1|t}\tilde{x}_{t+1|t}^\top] \\
 &= \mathbb{E}[(\alpha_t F_t(I - K_t \beta_t H_t)\tilde{x}_{t|t-1} + \xi_t) \\
 &\quad \times (\alpha_t F_t(I - K_t \beta_t H_t)\tilde{x}_{t|t-1} + \xi_t)^\top] \\
 &= [\alpha_t F_t(I - K_t \beta_t H_t)]\hat{P}_{t|t-1} \\
 &\quad \times [\alpha_t F_t(I - K_t \beta_t H_t)]^\top + \Delta P_{t+1|t}^{(1)} + Q_t.
 \end{aligned} \tag{13}$$

Here $\Delta P_{t+1|t}^{(1)}$ is defined as the covariance bias from expectation:

$$\begin{aligned}
 \Delta P_{t+1|t}^{(1)} &= \mathbb{E}\{[\alpha_t F_t(I - K_t \beta_t H_t)\tilde{x}_{t+1|t}] \\
 &\quad \times [\alpha_t F_t(I - K_t \beta_t H_t)\tilde{x}_{t+1|t}]^\top\} \\
 &\quad - [\alpha_t F_t(I - K_t \beta_t H_t)]\hat{P}_{t|t-1}[\alpha_t F_t(I - K_t \beta_t H_t)].
 \end{aligned} \tag{14}$$

Let $\Delta P_{t+1|t}^{(2)}$ be the difference between the real prediction error covariance $P_{t+1|t}$ and the estimated one $\hat{P}_{t+1|t}$ in (6):

$$\Delta P_{t+1|t}^{(2)} = P_{t+1|t} - \hat{P}_{t+1|t}. \tag{15}$$

Using (13)–(15), we can obtain (10a) and (11a) with $\Delta P_{t+1|t} = \Delta P_{t+1|t}^{(1)} - \Delta P_{t+1|t}^{(2)}$.

Similarly, we can derive the actual measurement error covariance from (9b).

$$\begin{aligned}
 P_{yy,t+1} &= \mathbb{E}[\tilde{y}_{t+1}\tilde{y}_{t+1}^\top] \\
 &= \mathbb{E}[(\beta_{t+1}H_{t+1}\tilde{x}_{t|t-1} + \zeta_t) \\
 &\quad \times (\beta_{t+1}H_{t+1}\tilde{x}_{t|t-1} + \zeta_t)^\top] \\
 &= (\beta_{t+1}H_{t+1})\hat{P}_{t+1|t}(\beta_{t+1}H_{t+1})^\top + \Delta P_{yy,t+1}^{(1)}, \\
 \Delta P_{yy,t+1}^{(1)} &= \mathbb{E}[(\beta_{t+1}H_{t+1})\tilde{x}_{t+1|t}\tilde{x}_{t+1|t}^\top(\beta_{t+1}H_{t+1})^\top] \\
 &\quad - (\beta_{t+1}H_{t+1})\hat{P}_{t+1|t}(\beta_{t+1}H_{t+1})^\top.
 \end{aligned} \tag{16}$$

Define $\Delta P_{yy,t+1}^{(2)} = P_{yy,t+1} - \hat{P}_{yy,t+1}$, $\Delta P_{yy,t+1} = \Delta P_{yy,t+1}^{(1)} - \Delta P_{yy,t+1}^{(2)}$, and we obtain (10b) (11b).

APPENDIX B PROOF OF PROPOSITION 2

Based on Assumption 2, we can derive the low bound of $\hat{P}_{t+1|t}$ from (7e):

$$\begin{aligned}
 \hat{P}_{t+1|t} &= \hat{P}_{t+1} + K_t \hat{P}_{xy,t+1}^\top \\
 &= \hat{P}_{t+1} + \hat{P}_{xy,t+1} \hat{P}_{yy,t+1}^{-1} \hat{P}_{xy,t+1}^\top \\
 &\geq \hat{P}_{t+1} \geq p_l I.
 \end{aligned}$$

Under Assumption 1, we can calculate the upper bound of $\hat{P}_{t+1|t}$ as (17).

$$\begin{aligned}
 \hat{P}_{t+1|t} &\leq \alpha_t F_t \hat{P}_{t+1} (\alpha_t F_t)^\top + Q_t + \Delta P_{t+1|t} \\
 &\leq (p_u \alpha_u^2 f_u^2 + q_u + \Delta p_u) I.
 \end{aligned} \tag{17}$$

According to Cauchy-Schwartz inequality, we can prove

$$\hat{P}_{xy,t+1} \leq \hat{P}_{t+1|t} \hat{P}_{yy,t+1}. \tag{18}$$

Substituting (18) in (7c) gives

$$\|K_t\| \leq \|\hat{P}_{t+1|t}\| = \sqrt{n}(p_u \alpha_u^2 f_u^2 + q_u + \Delta p_u) = K_u,$$

where n is the dimension of state x_t .

APPENDIX C PROOF OF THEOREM 2

For analyzing the boundedness of stochastic processes, the lemma is recalled:

Lemma 2. (Stability of stochastic process [38]) Assume that there is a stochastic process ψ_t with transform form $V(\psi_t)$, and real constants $\nu_l, \nu_u, \mu \geq 0$ and $0 \leq \lambda \leq 1$ such that $\forall k$

$$\nu_l \|\psi_t\|^2 \leq V(\psi_t) \leq \nu_u \|\psi_t\|^2, \quad (19a)$$

$$\mathbb{E}[V(\psi_t)|\psi_{t-1}] - V(\psi_{t-1}) \leq \mu - \lambda V(\psi_{t-1}). \quad (19b)$$

Then the process ψ_t is bounded in the mean square, i.e.,

$$\mathbb{E}\{\|\psi_t\|^2\} \leq \frac{\nu_u}{\nu_l} \mathbb{E}\{\|\psi_0\|^2\} (1-\lambda)^k + \frac{\mu}{\nu_l} \sum_{i=1}^{k-1} (1-\lambda)^i.$$

Define $V_{t+1}(\tilde{x}_{t+1|t}) = \tilde{x}_{t+1|t}^\top \hat{P}_{t+1|t}^{-1} \tilde{x}_{t+1|t}$ and according to Proposition 2, we have $\frac{\|\tilde{x}_{t+1|t}\|^2}{p_u \alpha_u^2 f_u^2 + q_u + \Delta p_u} \leq V_{t+1}(\tilde{x}_{t+1|t})^2 \leq \frac{\|\tilde{x}_{t+1|t}\|^2}{p_l}$, which fulfill the first condition (19a) of Lemma 2 if we set:

$$\nu_l = \frac{1}{p_u \alpha_u^2 f_u^2 + q_u + \Delta p_u}, \quad \nu_u = \frac{1}{p_l}. \quad (20)$$

Now let's consider :

$$\begin{aligned} & \mathbb{E}\{V_{t+1}(\tilde{x}_{t+1|t})|\tilde{x}_{t|t-1}\} \\ &= \tilde{x}_{t|t-1}^\top [\alpha_t F_t (I - K_t \beta_t H_t)]^\top \hat{P}_{t+1|t}^{-1} \\ & \times [\alpha_t F_t (I - K_t \beta_t H_t)] \tilde{x}_{t|t-1} \\ &+ \mathbb{E}\{(\zeta_t \alpha_t F_t K_t)^\top \hat{P}_{t+1|t}^{-1} \alpha_t F_t K_t \zeta_t | \tilde{x}_{t|t-1}\} \\ &+ \mathbb{E}\{\xi_t^\top \hat{P}_{t+1|t}^{-1} \xi_t | \tilde{x}_{t|t-1}\}. \end{aligned} \quad (21)$$

- Consider the first term of (21): $(\star) = [\alpha_t F_t (I - K_t \beta_t H_t)]^\top \hat{P}_{t+1|t}^{-1} [\alpha_t F_t (I - K_t \beta_t H_t)]$. Taking the inverse operation on it and substituting (10a) into it, the function becomes

$$\begin{aligned} (\star)^{-1} &= \hat{P}_{t|t-1} + [\alpha_t F_t (I - K_t \beta_t H_t)]^{-1} \hat{Q}_t \\ & \times [\alpha_t F_t (I - K_t \beta_t H_t)] \\ & \geq [1 + \frac{(q_l + \Delta p_l)I}{\hat{p}_u(\alpha_u f_u + \alpha_u f_u K_u \beta_u h_u)^2}] \hat{P}_{t|t-1}, \\ (\star) &\leq [1 + \frac{(q_l + \Delta p_l)I}{\hat{p}_u(\alpha_u f_u + \alpha_u f_u K_u \beta_u h_u)^2}]^{-1} \hat{P}_{t|t-1}^{-1}. \end{aligned} \quad (22)$$

Conveniently, define λ :

$$\lambda = 1 - [1 + \frac{(q_l + \Delta p_l)I}{\hat{p}_u(\alpha_u f_u + \alpha_u f_u K_u \beta_u h_u)^2}]^{-1}, \quad (23)$$

where obviously $0 \leq \lambda \leq 1$.

- Due to Assumption 2, the remainder of (21) can be derived as (24).

$$\begin{aligned} & \mathbb{E}\{(\zeta_t \alpha_t F_t K_t)^\top \hat{P}_{t+1|t}^{-1} \alpha_t F_t K_t \zeta_t | \tilde{x}_{t|t-1}\} \\ &+ \mathbb{E}\{\xi_t^\top \hat{P}_{t+1|t}^{-1} \xi_t | \tilde{x}_{t|t-1}\} \\ &\leq \mathbb{E}\left\{\frac{K_u^2 \alpha_u^2 f_u^2}{p_l} \text{tr}\{\zeta_t^T \zeta_t\} + \frac{1}{p_l} \text{tr}\{\xi_t^T \xi_t\}\right\} \\ &\leq \frac{K_u^2 \alpha_u^2 f_u^2 (r_u + \Delta r_u)}{p_l} m + \frac{q_u}{p_l} n = \mu. \end{aligned} \quad (24)$$

- Under the manipulation (21)–(24), we get that $\mathbb{E}\{V_{t+1}(\tilde{x}_{t+1|t})|\tilde{x}_{t|t-1}\} - V_t(\tilde{x}_{t|t-1}) \leq -\lambda V_t(\tilde{x}_{t|t-1}) + \mu$, which fulfill the condition (19b). Obviously, apply Lemma 2, and the estimation error $\tilde{x}_{t+1|t}$ is bounded in the mean square:

$$\mathbb{E}\{\|\tilde{x}_{t+1|t}\|^2\} \leq \frac{\nu_u}{\nu_l} \mathbb{E}\{\|x_0\|^2\} (1-\lambda)^t + \frac{\mu}{\nu_l} \sum_{i=1}^{t-1} (1-\lambda)^i.$$

From (7d), we can yield \tilde{x}_{t+1} is similarly bounded:

$$\begin{aligned} \mathbb{E}\{\|\tilde{x}_{t+1}\|^2\} &\leq (1 + f_u K_u \beta_u h_u) \mathbb{E}\{\|\tilde{x}_{t+1|t}\|^2\} \\ &+ r_u + \Delta r_u + \Delta p_{yy,u}. \end{aligned} \quad (25)$$

Analyze (20) (23) (24) and (25), we can obtain how $r_u, \Delta r_u$ and $\mathbb{E}\{\|x_0\|^2\}$ affect the bound.

$$\mathbb{E}\{\|\tilde{x}_{t+1}\|^2\} \leq C_1 \mathbb{E}\{\|x_0\|^2\} + C_2 (r_u + \Delta r_u) + C_3,$$

where C_1, C_2, C_3 are constants as follows:

$$\begin{aligned} C_1 &= \frac{\nu_u}{\nu_l} (1-\lambda) (1 + f_u K_u \beta_u h_u), \\ C_2 &= 1 + (1 + f_u K_u \beta_u h_u) \frac{m K_u^2 \alpha_u^2 f_u^2}{\nu_l p_l (1-\lambda)}, \\ C_3 &= (1 + f_u K_u \beta_u h_u) \frac{n q_u}{\nu_l p_l (1-\lambda)} + \Delta p_{yy,u}. \end{aligned} \quad (26)$$

The relationship between the bound of estimation error and $r_u, \Delta r_u$, and $\mathbb{E}\{x_0\}^2$ is evidently linearly positive. It is notable that the analyzed bound isn't the supremacy.

REFERENCES

- [1] W. Luo, J. Xing, A. Milan, X. Zhang, W. Liu, and T.-K. Kim, "Multiple object tracking: A literature review," *Artificial intelligence*, vol. 293, p. 103448, 2021.
- [2] S. Liu, S. Huang, X. Xu, J. Lloret, and K. Muhammad, "Efficient visual tracking based on fuzzy inference for intelligent transportation systems," *IEEE Transactions on Intelligent Transportation Systems*, 2023.
- [3] S. E. Li, *Reinforcement learning for sequential decision and optimal control*. Springer, 2023.
- [4] X. Weng and K. Kitani, "A baseline for 3d multi-object tracking," *arXiv preprint arXiv:1907.03961*, vol. 1, no. 2, p. 6, 2019.
- [5] G. Guo and S. Zhao, "3d multi-object tracking with adaptive cubature kalman filter for autonomous driving," *IEEE Transactions on Intelligent Vehicles*, vol. 8, no. 1, pp. 512–519, 2022.
- [6] X. Li, T. Xie, D. Liu, J. Gao, K. Dai, Z. Jiang, L. Zhao, and K. Wang, "Poly-mot: A polyhedral framework for 3d multi-object tracking," in *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 9391–9398, IEEE, 2023.
- [7] L. Wang, X. Zhang, W. Qin, X. Li, J. Gao, L. Yang, Z. Li, J. Li, L. Zhu, H. Wang, et al., "Camo-mot: Combined appearance-motion optimization for 3d multi-object tracking with camera-lidar fusion," *IEEE Transactions on Intelligent Transportation Systems*, 2023.
- [8] X. Li, D. Liu, L. Zhao, Y. Wu, X. Wu, and J. Gao, "Fast-poly: A fast polyhedral framework for 3d multi-object tracking," *arXiv preprint arXiv:2403.13443*, 2024.

- [9] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, "Vision meets robotics: The kitti dataset," *The International Journal of Robotics Research*, vol. 32, no. 11, pp. 1231–1237, 2013.
- [10] S. Shi, X. Wang, and H. Li, "Pointcnn: 3d object proposal generation and detection from point cloud," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 770–779, 2019.
- [11] S. K. Pal, A. Pramanik, J. Maiti, and P. Mitra, "Deep learning in multi-object detection and tracking: state of the art," *Applied Intelligence*, vol. 51, pp. 6400–6429, 2021.
- [12] L. Rakai, H. Song, S. Sun, W. Zhang, and Y. Yang, "Data association in multiple object tracking: A survey of recent techniques," *Expert Systems with Applications*, vol. 192, p. 116300, 2022.
- [13] N. V. Dung, N. L. Trung, K. Abed-Meraim, et al., "Robust subspace tracking with missing data and outliers: Novel algorithm with convergence guarantee," *IEEE Transactions on Signal Processing*, vol. 69, pp. 2070–2085, 2021.
- [14] S. Pang, D. Morris, and H. Radha, "3d multi-object tracking using random finite set-based multiple measurement models filtering (rfs-m 3) for autonomous vehicles," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 13701–13707, IEEE, 2021.
- [15] W. Cao, C. Liu, Z. Lan, Y. Piao, and S. E. Li, "Generalized moving horizon estimation for nonlinear systems with robustness to measurement outliers," in *2023 American Control Conference (ACC)*, pp. 1614–1621, IEEE, 2023.
- [16] T. Omeragić and J. Velagić, "Tracking of moving objects based on extended kalman filter," in *2020 International Symposium ELMAR*, pp. 137–140, IEEE, 2020.
- [17] M. Liu, J. Niu, and Y. Liu, "Ukf-mot: An unscented kalman filter-based 3d multi-object tracker," *CAAI Transactions on Intelligence Technology*, 2024.
- [18] Z. Li, W. Wang, H. Li, E. Xie, C. Sima, T. Lu, Y. Qiao, and J. Dai, "Bevformer: Learning bird's-eye-view representation from multi-camera images via spatiotemporal transformers," in *European conference on computer vision*, pp. 1–18, Springer, 2022.
- [19] Y. Qian, X. Wang, H. Zhuang, C. Wang, and M. Yang, "3d vehicle detection enhancement using tracking feedback in sparse point clouds environments," *IEEE Open Journal of Intelligent Transportation Systems*, 2023.
- [20] Y. Li, J. Zhu, S. C. Hoi, W. Song, Z. Wang, and H. Liu, "Robust estimation of similarity transformation for visual object tracking," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 33, pp. 8666–8673, 2019.
- [21] 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 multi-modal dataset for autonomous driving," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 11621–11631, 2020.
- [22] H. W. Kuhn, "The hungarian method for the assignment problem," *Naval research logistics quarterly*, vol. 2, no. 1-2, pp. 83–97, 1955.
- [23] S. E. Li, G. Li, J. Yu, C. Liu, B. Cheng, J. Wang, and K. Li, "Kalman filter-based tracking of moving objects using linear ultrasonic sensor array for road vehicles," *Mechanical Systems and Signal Processing*, vol. 98, pp. 173–189, 2018.
- [24] R. Schubert, E. Richter, and G. Wanielik, "Comparison and evaluation of advanced motion models for vehicle tracking," in *2008 11th International Conference on Information Fusion*, pp. 1–6, 2008.
- [25] S. J. Julier and J. K. Uhlmann, "Unscented filtering and nonlinear estimation," *Proceedings of the IEEE*, vol. 92, no. 3, pp. 401–422, 2004.
- [26] W. Cao, S. Liu, C. Liu, Z. He, S. S.-T. Yau, and S. E. Li, "Convolutional bayesian filtering," *arXiv preprint arXiv:2404.00481*, 2024.
- [27] S.-C. Pei and J.-J. Ding, "Closed-form discrete fractional and affine fourier transforms," *IEEE transactions on signal processing*, vol. 48, no. 5, pp. 1338–1353, 2000.
- [28] K. Xiong, H. Zhang, and C. Chan, "Performance evaluation of ukf-based nonlinear filtering," *Automatica*, vol. 42, no. 2, pp. 261–270, 2006.
- [29] L. Li and Y. Xia, "Stochastic stability of the unscented kalman filter with intermittent observations," *Automatica*, vol. 48, no. 5, pp. 978–981, 2012.
- [30] K. Reif, S. Gunther, E. Yaz, and R. Unbehauen, "Stochastic stability of the discrete-time extended kalman filter," *IEEE Transactions on Automatic control*, vol. 44, no. 4, pp. 714–728, 1999.
- [31] B. Zhu, Z. Jiang, X. Zhou, Z. Li, and G. Yu, "Class-balanced grouping and sampling for point cloud 3d object detection," *arXiv preprint arXiv:1908.09492*, 2019.
- [32] Z. Bing, C. Lubin, J. Xu, et al., "Huber-based adaptive unscented kalman filter with non-gaussian measurement noise [j]," *Circuits Systems and Signal Processing*, vol. 37, no. 12, pp. 1–21, 2018.
- [33] B. Ge, H. Zhang, L. Jiang, Z. Li, and M. M. Butt, "Adaptive unscented kalman filter for target tracking with unknown time-varying noise covariance," *Sensors*, vol. 19, no. 6, p. 1371, 2019.
- [34] J. Havlík and O. Straka, "Performance evaluation of iterated extended kalman filter with variable step-length," in *Journal of Physics: Conference Series*, vol. 659, p. 012022, IOP Publishing, 2015.
- [35] Z. Qiu, S. Wang, P. Hu, and L. Guo, "Outlier-robust extended kalman filtering for bioinspired integrated navigation system," *IEEE Transactions on Automation Science and Engineering*, 2023.
- [36] Z. Qiu and L. Guo, "Improved cubature kalman filter for spacecraft attitude estimation," *IEEE Transactions on Instrumentation and Measurement*, vol. 70, pp. 1–13, 2020.
- [37] K. Bernardin and R. Stiefelhagen, "Evaluating multiple object tracking performance: the clear mot metrics," *EURASIP Journal on Image and Video Processing*, vol. 2008, pp. 1–10, 2008.
- [38] T.-J. Tarn and Y. Rasis, "Observers for nonlinear stochastic systems," *IEEE Transactions on Automatic Control*, vol. 21, no. 4, pp. 441–448, 1976.

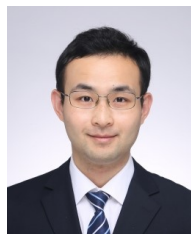


Shiqi Liu received his B.E. degree in the School of Vehicle and Mobility from Tsinghua University, Beijing, China, in 2023. He is currently a Ph.D. candidate in the School of Vehicle and Mobility at Tsinghua University, Beijing, China. His research interests include optimal filtering, Bayesian inference, and reinforcement learning.



Wenhan Cao received his B.E. degree in the School of Electrical Engineering from Beijing Jiaotong University, Beijing, China, in 2019.

He is currently a Ph.D. candidate in the School of Vehicle and Mobility, Tsinghua University, Beijing, China. His research interests include optimal filtering and reinforcement learning. He was a finalist for the Best Student Paper Award at the 2021 IFAC MECC.



Chang Liu (Member, IEEE) received the B.S. degrees in Electronic Information Science and Technology and in Mathematics and Applied Mathematics (double degree) from the Peking University, China, in 2011, and the M.S. degrees in Mechanical Engineering and in Computer Science, and the Ph.D. degree in Mechanical Engineering from the University of California, Berkeley, USA, in 2014, 2015, and 2017, respectively.

He is currently an Assistant Professor with the Department of Advanced Manufacturing and Robotics, College of Engineering, Peking University. From 2017 to 2020, he was a Postdoctoral Associate with the Cornell University, USA. He has also worked for Ford Motor Company and NVIDIA Corporation on autonomous vehicles. His research interests include robot motion planning, active sensing, and multi-robot collaboration.



Tianyi Zhang received his B.E. degree in Automation Science and Electrical Engineering from Beihang University, Beijing, China, in 2024. He will soon begin his Ph.D. studies in the School of Vehicle and Mobility at Tsinghua University, Beijing, China. His research interests include optimal state estimation, Bayesian inference, and reinforcement learning.



Shengbo Eben Li (Senior Member, IEEE) received his M.S. and Ph.D. degrees from Tsinghua University in 2006 and 2009. Before joining Tsinghua University, he has worked at Stanford University, University of Michigan, and UC Berkeley. His active research interests include intelligent vehicles and driver assistance, deep reinforcement learning, optimal control and estimation, etc. He is the author of over 190 peer-reviewed journal/conference papers, and co-inventor of over 40 patents. Dr. Li has received over 20 prestigious awards, including Youth

Sci. & Tech Award of Ministry of Education (annually 10 receivers in China), Natural Science Award of Chinese Association of Automation (First level), National Award for Progress in Sci & Tech of China, and best (student) paper awards of IET ITS, IEEE ITS, IEEE ICUS, CVCI, etc. He also serves as Board of Governor of IEEE ITS Society, Senior AE of IEEE OJ ITS, and AEs of IEEE ITSM, IEEE TITS, IEEE TIV, IEEE TNNLS, etc.