# MSight: An Edge-Cloud Infrastructure-based Perception System for Connected Automated Vehicles

Rusheng Zhang*, Depu Meng*, *Member, IEEE*, Shengyin Shen, Zhengxia Zou,
Houqiang Li, *Fellow, IEEE* , Henry X. Liu, *Member, IEEE*

*Abstract*—As vehicular communication and networking technologies continue to advance, infrastructure-based roadside perception emerges as a pivotal tool for connected automated vehicle (CAV) applications. Due to their elevated positioning, roadside sensors, including cameras and lidars, often enjoy unobstructed views with diminished object occlusion. This provides them a distinct advantage over onboard perception, enabling more robust and accurate detection of road objects. This paper presents MSight, a cutting-edge roadside perception system specifically designed for CAVs. MSight offers real-time vehicle detection, localization, tracking, and short-term trajectory prediction. Evaluations underscore the system's capability to uphold lane-level accuracy with minimal latency, revealing a range of potential applications to enhance CAV safety and efficiency. Presently, MSight operates 24/7 at a two-lane roundabout in the City of Ann Arbor, Michigan.

*Index Terms*—Roadside perception, vehicle-to-infrastructure communications, cooperative perception, connected automated vehicle, autonomous vehicle
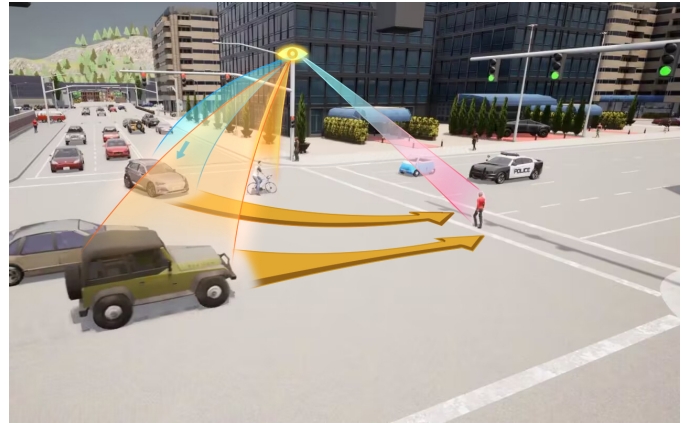
Fig. 1: An illustrative figure of cooperative perception. The roadside sensor detects vehicles on the road, and the perception results are forwarded to connected vehicles via roadside radio.

## I. INTRODUCTION

With the rapid development in vehicular communication and networking technologies, infrastructure-based roadside vehicle detection has become viable for connected automated vehicle (CAV) applications. Roadside perception systems can be more accurate, reliable, and with better coverage than vehicle onboard perception systems due to smaller background variations and elevated positions. Connected vehicles can receive the roadside perception results via vehicle-to-infrastructure (V2I) and vehicle-to-everything (V2X) communication with short latency to improve the perception quality cooperatively, as shown in Figure 1.

While many current automated vehicle solutions rely on onboard perception systems, their limitations become increasingly evident as the demands for more complex perception tasks grow. Occlusion stands out as a significant drawback of onboard perception, potentially leading to safety-critical situations [1]. Computational power onboard also constrains the application of resource-intensive perception algorithms. In contrast, roadside perception systems can employ high-performance computing hardware, facilitating more sophisticated algorithms and minimizing inference latency. Although they bear similarities to surveillance and monitoring systems, roadside perception systems for connected and automated vehicles distinctly demand high localization and tracking accuracy coupled with low latency.

Roadside perception systems also have significant impacts in the realm of vehicular communication. Through V2V communications, vehicles broadcast Basic Safety Messages (BSMs) as dictated by the SAE 2735 protocol [2]. These messages are broadcast either via the 802.11p channel or through C-V2X technology [3], [4]. These broadcasts relay crucial safety data, such as location, heading, brake status, and more, which neighboring vehicles can utilize. Yet, the uptake of this technology is sluggish, largely due to a classic

R. Zhang is with the Department of Civil and Environmental Engineering, University of Michigan, Ann Arbor, MI, 48109, USA (email: rushengz@umich.edu)

D. Meng is with the Department of Civil and Environmental Engineering, University of Michigan, Ann Arbor, MI, 48109, USA (email: depum@umich.edu)

S. Shen is with the University of Michigan Transportation Research Institude, 2901 Baxer Rd, Ann Arbor, MI, 48109, USA. (email: shengyin@umich.edu)

Z. Zou is with the Department of Guidance, Navigation and Control, School of Astronautics, Beihang University, Beijing 100191, China. (email: zhengxiazou@buaa.edu.cn)

Houqiang Li is with the CAS Key Laboratory of Technology in Geo-Spatial Information Processing and Application System, Department of Electronic Engineering and Information Science, University of Science and Technology of China, Hefei, China (email: lihq@ustc.edu.cn)

H. X. Liu is with the Department of Civil and Environmental Engineering, University of Michigan, Ann Arbor, MI, 48109, USA. H. X. Liu is also with Mcity, University of Michigan, Ann Arbor, MI, 48109, USA (email: henryliu@umich.edu)

*Equal contribution.

'chicken-and-egg' dilemma. Its sparse adoption offers minimal advantages to early adopters, leading to hesitancy among manufacturers, which perpetuates the cycle of slow adoption. To overcome this issue, proxy solutions via V2I communications are proposed to boost penetration rate at certain locations. Roadside radios, termed as Roadside Units (RSUs), encode these detection results, disseminating them to all CAVs in the vicinity. This creates an effect synonymous with high adoption within chosen regions. An initial method involves encoding perception data to resemble broadcasts from vehicles, referred to as Proxy-BSM in the USA or Proxy-CAM (cooperative awareness message) in Europe and Japan [5]. More recently, the SAE J3224 standard has been put forward, offering a structured approach using Sensor Data Sharing Messages (SDSMs) [6].

Currently, roadside perception systems for cooperative driving is still in its infancy and many reported results are still in lab or early testing stage. The existing works are preliminary and far from real-world deployment. In this paper, we introduce MSight, a full-stack end-to-end roadside perception system. MSight stands out from existing solutions by aiming to achieve high accuracy and low latency that meets the requirements of CAV applications.

There are several aspects of MSight that make it an imperative step towards a production-ready cooperative perception system. First of all, MSight is a full-stack end-to-end cooperative perception system that contains not only a state-of-the-art perception algorithm but also a complete communication pipeline with CAVs and the cloud data center. Furthermore, the system has been deployed and operated 24/7 on the roadside in a production environment. Finally, this system has been extensively tested with CAVs, and the field test results are presented in this paper. It is anticipated that MSight will accelerate the deployment of CAVs and speed up the adoption of V2X communications technologies. Despite its primary focus on CAV applications, the edge-cloud architecture of MSight makes it possible to extend its functionalities in a wide range of different applications.

## II. RELATED WORKS

Sensor-based roadside surveillance systems date back to 1986. Initially, these systems were used to detect abnormal vehicle behavior [7]. The development of roadside surveillance/perception systems has been rapid since then. Such systems typically use one or more cameras mounted at an elevated position on the roadside to detect and track moving objects as well as detect traffic conflicts [8]. To detect road objects, different methods have been used, to name a few, background subtraction [9], frame difference [10], synthesizing training data with AR and GAN [11], feature-based detection [12], KanadeLucas-Tomasi tracking [13], cascading classifiers [14] and many more [7].

In recent years, with the development of deep learning (DL) techniques in computer vision, real-time high-quality vision-based perception algorithms are proposed and widely applied into ADAS systems as well as automated vehicles. Though transformer-based object detectors achieved state-of-the-art detection accuracy [15]–[20], the YOLO [21]–[25] series object

detectors are known for lightweight and fast processing and suitable for CAV applications. Vehicle trajectories are valuable data for transportation research. Visual tracking algorithms are designed to obtain the object trajectories from consecutive video frames. There are two classes of object tracking. single object tracking (SoT) [26]–[31] and multiple object tracking (MoT) [32]–[37]. MoT methods can be deployed in roadside perceptions for object trajectory tracking.

These new DLs are also adopted in roadside vehicle perception and cooperative driving. For instance, [38] introduces a method that detects vehicles from roadside based on YOLOv5 detector and finds anomaly behavior using decision tree. [39] proposes a modified version of YOLOv5, that optimizes the performance for roadside perception tasks. In general, DL-based methodology for vehicle detection is at the initial stage, but of significant potential, thereby attracting increasing attention.

In addition to the aforementioned systems using regular cameras, a wide range of other sensors have also been explored to overcome the limitations of regular cameras. For example, systems with other types of cameras are investigated, for instance, fish-eye cameras to enhance per-camera coverage [40], and thermal cameras to increase the robustness at night and in different weathers [41], [42]. Recently, roadside lidars are explored in several researches for roadside vehicle detection [43]. These methods include traditional point cloud detection methods based on background subtraction and point clustering [44]–[46], as well as DL-based methods [47], [48]. [49], [50] proposes lidar-based vehicle detection methods that are robust in adverse weather, especially in snow and heavy rain.

Furthermore, methods using multiple sensors and fusion strategies have been proposed for roadside perception. In [51], a roadside perception system with lidar and dual cameras is prototyped and demonstrated with CAVs self-driving using the roadside perception information only. In [52], [53], multiple fish-eye cameras and thermal cameras are utilized for vehicle detection. [54] provide a sensor fusion method that synchronizes detection results with millimeter-wave radar and camera detection. For some CAV applications, the ability of predicting the objects' future positions is also required. Trajectory prediction methods often use objects' historical trajectories with semantic map information to predict objects' future trajectories [55]–[58] or directly use raw sensor data to predict future trajectories [59], [60].

On the other hand, V2X communications are proposed to transmit roadside perception results to CAVs in a real-time fashion. [61] provides a design using a V2X communication system for cooperative perception, and reports experimental results conducted by 802.11p radios and measures the transmission latency of cooperative awareness messages (CAM). Their research has found that the transmission latency of CAMs is within 20ms, but the overall latency can be as large as 200ms considering the idle time between periodical transmissions. In [62], a vehicle-roadside cooperative perception system is proposed with a high-level fusion strategy. In [63], the authors report the design of a communication system that broadcasts proxy-CAM to CVs and some initial field-tests results on communication reception rate and latency.
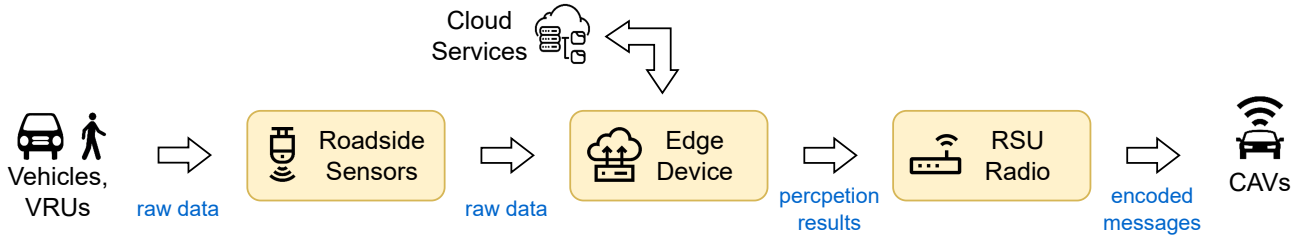
Fig. 2: An illustration of the roadside portion of MSight system

The work is continued in [64], [65], a roadside perception unit is implemented with open-source autonomous driving software, inter-connected with RSU and sensors with high-speed networks. Thorough latency analysis is carried out using OMNET++ and field experiments. [66] systematically designs scalable communication network structures for infrastructure-vehicle cooperative driving tasks.

All aforementioned works concentrate on distinct aspects of a roadside perception system, without systematic integration, testing, and field implementation. In this work, we endeavor to orchestrate a comprehensive framework, encompassing object detection, localization, tracking, prediction, vehicle list encoding, message forwarding, and vehicle-side decoding. This integrative approach is a pivotal advancement towards achieving scalable, production-grade deployments in real-world scenarios. Moreover, the evaluation methods used in this paper focus on end-to-end perception performance for CAV applications, offering a more extensive exploration compared to existing works that primarily resolve singular issues or components. Consequently, this research serves not only as an academic advancement but also as a meticulously crafted template, paving the way for pragmatic deployments in the field of roadside perception.

## III. SYSTEM DESIGN

### A. The Edge-Cloud System Design

As mentioned above, MSight is a specialized, real-time roadside perception system, predominantly devised for CAVs. It bifurcates into two critical components: the roadside portion and the cloud portion. The former revolves around an edge device, serving as the nexus that orchestrates all roadside sensors, executes delay-sensitive applications, and fosters communication with CAVs through RSUs. Concurrently, the cloud component is engineered to be highly scalable, capable of assimilating expansive data streams. Data are archived in cloud storage and disseminated to various services via a publish/subscribe mechanism within the cloud. Together, these two portions construct an intricate edge-cloud system architecture, endowed with the versatility to accommodate a diverse array of downstream applications.

### B. The Roadside Portion

The principal objective of the roadside portion is to execute the perception algorithm, detailed in section IV, and to convey the resulting perceptions to CAVs with minimal latency, as illustrated in Figure 2. Vehicles and Vulnerable Road Users

(VRUs) are identified through roadside sensors and the raw data they yield (such as raw images) are channeled to the edge device. This device processes the raw data using the perception algorithm and transmits the encoded perception results to CAVs via the RSU radio. This entire perception loop is localized to the roadside to guarantee minimal latency in the communication network. Simultaneously, the edge device relays selected data to the cloud, catering to the needs of other applications and services residing therein.

Figure 3 displays the array of devices installed at the roadside. While MSight is engineered to accommodate a wide range of sensors and radios, the devices discussed in this paper pertain specifically to our installation described in section III-D. Generally, the roadside component is subdivided into three primary elements: an edge device that operates as the nucleus, sensors that acquire raw data, and a radio that communicates with CAVs. Figure 3a illustrates the Nuvo edge device situated in the traffic cabinet near the roundabout, interconnected with sensors and RSU radio within the same network. GridSmart fisheye cameras are selected (as shown in Figure 3b), for their superior field of view compared to conventional pin-hole cameras. These cameras are strategically positioned at each quadrant (north-east, north-west, south-east, south-west) of the roundabout. Figure 3c presents the Cohda MK5 RSU radio, which is commissioned to periodically transmit messages encoded by the edge device.

### C. The Cloud Portion

Figure 4 illustrates the architectural design of MSight's cloud segment, constructed for extensibility and versatility. It incorporates a gateway to establish an API endpoint for the edge device, facilitating the streaming of roadside data. Once data reaches the gateway, it is relayed to a publication subscription (pub/sub) service, operating as a message dispatcher, distributing data to all its subscribers. This architecture ensures the decoupling of message producers and consumers, enhanc-



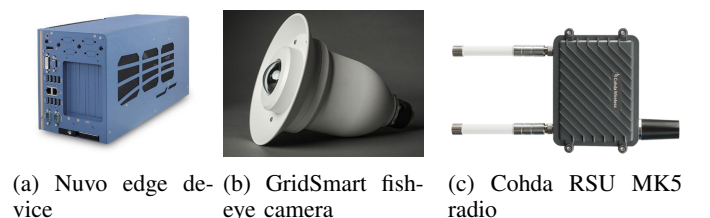(a) Nuvo edge device  (b) GridSmart fisheye camera  (c) Cohda RSU MK5 radio
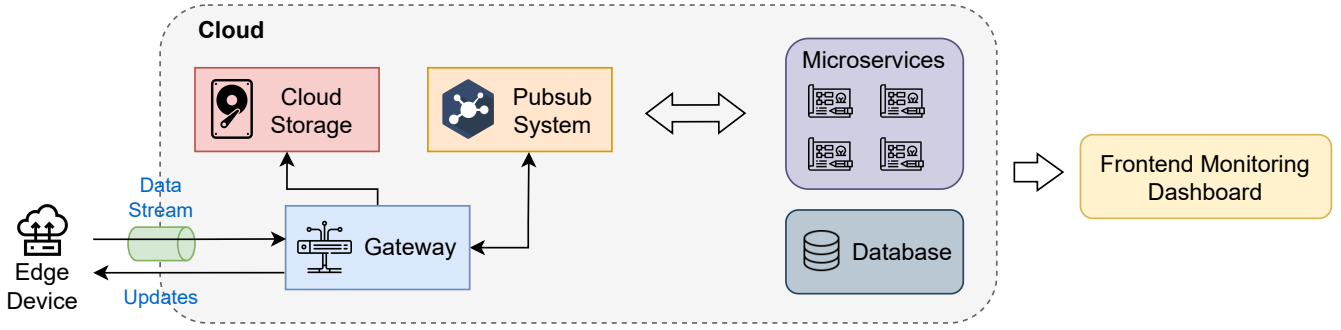
Fig. 3: Nuvo edge device

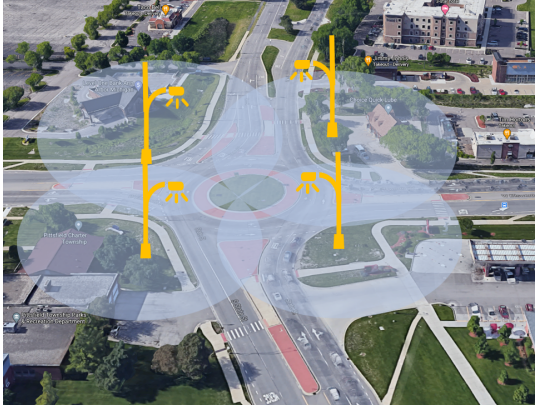Fig. 4: An illustration of the cloud-side portion of the MSight system



Fig. 5: The four corners of the roundabout installing the sensors, and the illustration of their coverage

ing the system's robustness, scalability, and error resilience. New services can be seamlessly integrated to consume the data stream without altering existing cloud components. The pub/sub system to forward messages with minimal latency, enabling near real-time services. This feature is particularly beneficial for third parties requiring instantaneous traffic data.

Each application within the cloud is constructed as a microservice, ensuring scalability and ease of maintenance. These services, using the aforementioned pub/sub system for data communication, are completely decoupled, enabling the streamlined development of additional services and affirming MSight as a highly extendable system.

Moreover, a distinct dedicated link is established between the gateway and cloud storage, allowing the direct storage of data streamed from roadside edge devices. This link is deliberately isolated from the pub/sub system and all other microservices residing on the cloud, safeguarding data security by design. This structured and secure approach ensures the integrity and reliability of MSight's cloud infrastructure, accommodating a range of applications and services.

### D. Deployment

This paper focuses on a particular deployment of MSight at an intersection between Ellsworth Road and State Street, Ann Arbor, Michigan, USA. Camera sensors are mounted on the light pole at the roadside of the roundabout. For this particular deployment, GridSmart fisheye cameras are selected due to

their large coverage area. A Cohda RSU radio is also mounted on the light pole. There are four fisheye cameras mounted at the four corners of the roundabout, as delineated in Figure 5. Additionally, a provisional coverage area for each sensor is represented in the figure, with each sensor monitoring a quarter of the roundabout.

## IV. PERCEPTION ALGORITHM

The perception algorithm aims at detecting, tracking, and predicting objects such as cars, trucks, buses, and motorcycles at the roundabout in real-time. The proposed framework consists of 6 components: camera calibration, image alignment, object detection, object localization, object tracking, and future trajectory prediction. As shown in Figure 6, the components are assembled in a sequential manner. Firstly, we calibrate the sensors to get the homography transformation matrices of sensors for image pixel location to real-world coordinate projection. Since the roadside camera view might slightly change over time, the calibration might be inaccurate. We align the input image to a standard camera view to guarantee the calibration accuracy. A YOLOX [25]-based lightweight object detector is applied for 2D object detection. For object localization, we first project the object pixel locations to real-world coordinates, then we fuse the object coordinates from multiple cameras together to form a unified object location result. A SORT [37]-based object tracking algorithm is used to extract object trajectories. With the historical trajectories, a transformer-based trajectory prediction module is applied to predict the object's future trajectories.

### A. Camera Calibration

Traditional checkerboard calibration methods are inefficient in this case, as they require field operations, specialized devices, and skilled operators on site. Therefore, we use a landmark-based calibration method. Figure 7 shows a set of landmark points selected for the calibration. Landmarks are identified both in the fisheye images produced by our fisheye camera and in satellite images in which the latitude and longitude of the points are known.

A calibration method is then developed using these landmark pairs from the two images. Figure 7 illustrates the calibration method. Assume that the camera lens follows a generic radially symmetric model. $r(\theta) = k_1\theta + k_2\theta^3 + k_3\theta^5 + \dots$
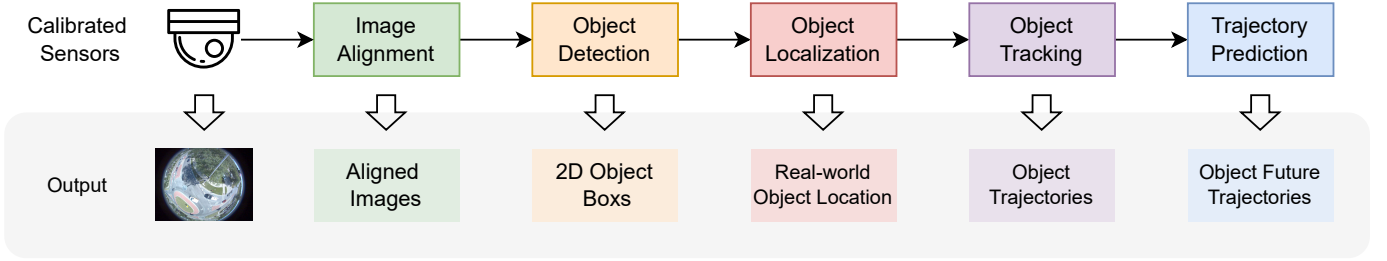
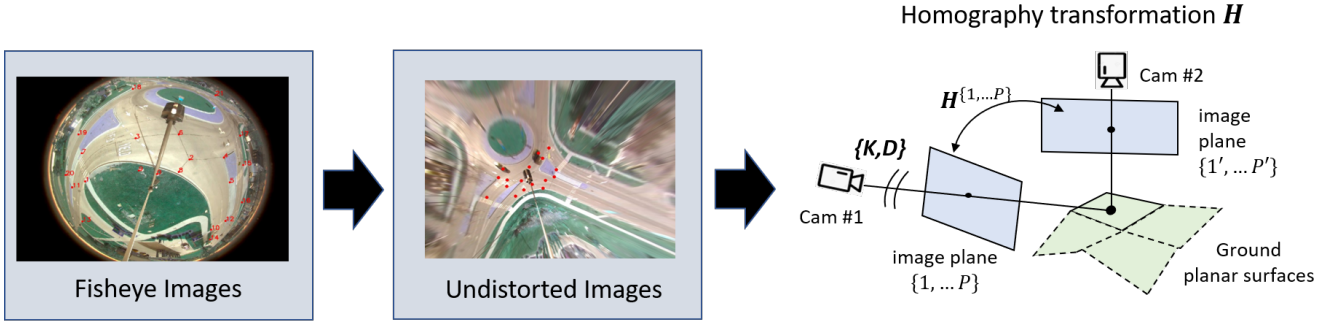Fig. 6: An illustration of overall proposed roadside perception framework



Fig. 7: The landmark-based calibration method for fisheye camera
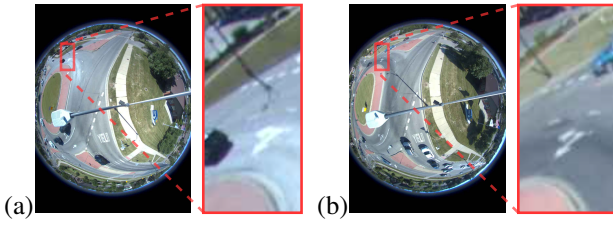


Fig. 8: An illustration of camera view change within the same day at (a) 10am and (b) 6pm
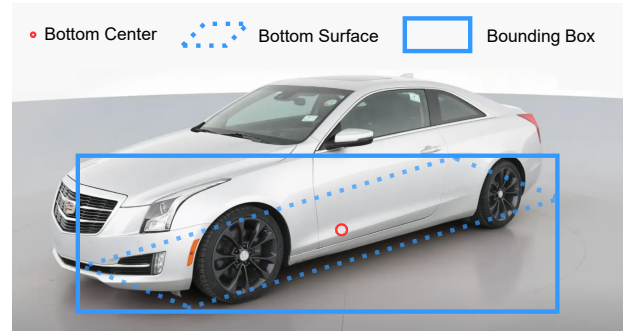


Fig. 9: An illustration of the bounding box annotation. We use the minimum bounding rectangle of object bottom surface as object bounding box for object detection.

where $\theta$ is the angle between the principle axis and the incoming ray, $r(\theta)$ is the distance between the correspondence point to the principle point [67]. An estimation of the intrinsic matrix of the camera and its distortion coefficients can be made according to [68]. Then, a homography transformation can be found between the undistorted image and the world latitude and longitude with least square regression and RANSAC consensus between the two groups of landmark sets. This transformation can be applied to map any point in the fisheye image to the world latitude and longitude.

### B. Image Alignment

The roadside cameras are installed on four poles at four corners of the roundabout. Since the roadside cameras will serve for an extended period of time once installed, both elastic deformation of the poles caused by temperature change, or metal creep of the poles could happen. This will lead to a change of camera view. As shown in Figure 8, The camera view at 10am and 6pm are different even on the same day. In the zoomed-in region, one can see that the arrow landmark at 6pm is higher than at 10am. These view changes, if not compensated, will lead to calibration errors.

To compensate for the view changes, we apply a run-time image alignment module. The module contains two parts: geometric transformation estimation and image wrapping. For geometric transformation estimation, given standard image $\mathbf{I}_s$ and input image $\mathbf{I}_i$, we first convert them to grayscale images $\hat{\mathbf{I}}_s$ and $\hat{\mathbf{I}}_i$. Then we use Enhanced Correlation Coefficient [69] to compute the geometric transformation $\mathbf{T}$ between two images.

$$\mathbf{T} = \arg\max_{\mathbf{W}} \mathrm{ECC}(\hat{\mathbf{I}}_i(x,y), \hat{\mathbf{I}}_s(x',y')) \tag{1}$$

Here

$$\begin{bmatrix} x' \\ y' \end{bmatrix} = \mathbf{W} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} \tag{2}$$
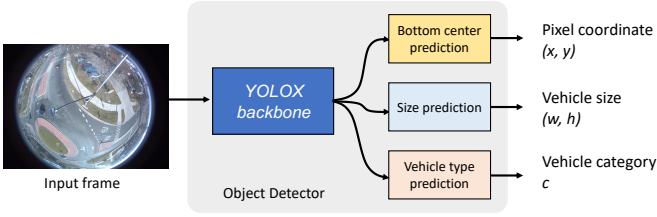
Fig. 10: An illustration of the object detection pipeline

With the geometric transformation $\mathbf{T}$, we can warp the input image to the same view as the standard image:

$$\tilde{\mathbf{I}}_i = \mathrm{WarpPerspective}(\mathbf{I}_i, \mathbf{T}) \tag{3}$$

Then we feed the wrapped image $\tilde{\mathbf{I}}_i$ into object detector for object detection.

### C. Object Detection

For object detection, we designed three object categories: car, truck/bus/trailer, and motorcyclist/cyclist/pedestrian. For object bounding boxes, instead of using the object 2d bounding boxes, we use the minimum bounding rectangle (MBR) of objects' bottom surface as illustrated in Figure 9. In this way, the center point of the bounding box will be the bottom center of the object.

*a) Network:* A standard YOLOX-nano [25] model is used for object detection. The YOLOX-nano is a lightweight single-stage object detector that contains $0.91\mathrm{M}$ parameters and $2.56\mathrm{G}$ FLOPs with input image size $640 \times 640$. The YOLOX-nano model contains a DarkNet backbone, a feature pyramid network (FPN) layer, and YOLOX detection head. The backbone extract deep feature from the input image, and the FPN layer aggregates multi-scale features output by the backbone. The multi-scale features are fed into the YOLOX detection head for classification and bounding box regression. Fig 10 shows an illustration of the object detection pipeline.

*b) Training:* For the object detector training, we use $5,000$ images in total as the training set, and the four camera views share the the same network. For data augmentation, the translation, sheer, rotation, and hsv augmentations are adopted. The model is trained for $150$ epochs with an initial learning rate $5e - 5$. The learning rate is dropped by a factor of 10 at epoch 100. The mini-batch size used for training is 8.

### D. Object Localization

Once we obtained the object 2D location in the image, we need to obtain the real-world object coordinates for downstream applications. In the object localization module, we first map the object 2D location to real-world latitude and longitude coordinates, then we fuse the object location information for multiple cameras together to get a merged object real-world locations across different camera views.

*a) Image to real-world mapping:* For image position $(x_p, y_p)$, we first obtain its undistorted position $(\hat{x}_p, \hat{y}_p)$ with the intrinsic parameter $K$ and distortion parameter $D$ of the fisheye camera. Then we apply the homography transformation $\mathbf{H}$ obtained by camera calibration to the undistorted image to get the real-world coordinate $(x_r, y_r)$.

*b) Multi-camera Fusion:* Since we installed four fisheye cameras at the four corners (Northeast, Northwest, Southeast, Southwest) at the roundabout, the detection results of the four cameras needs to be fused together to obtain a merged detection result of the whole roundabout. We split the roundabout into four parts (NE, NW, SE, SW) corresponding to the four camera locations. We assign each part as the region of interest of each camera. We only select the detected objects within the region of interest of each camera and put the detected objects together to obtain the fused detection results.

### E. Object Tracking

The object tracking module is built based on SORT [37], an online object tracking algorithm. The tracking module consists a Kalman Filter [70] for state estimation and a Hungarian Algorithm [71] for the association between existing targets and detections. Compared to the SORT algorithm, we use the predicted future locations of each target provided by our transformer-based trajectory prediction module instead of the Kalman Filter for association for better tracking performance. Following SORT, the state of each target is defined as

$$\mathbf{s} = [x_c, y_c, s, r, v_x, v_y, v_s, v_r]^\top \tag{4}$$

Here, $x_c$ and $y_c$ are the real-world latitude and longitude coordinate of the target, $s$ and $r$ are the scale and the aspect ratio of the target's bounding box respectively. $v_x, v_y, v_s, v_r$ are the derivatives of $x, y, s, r$.

To assign detections to existing targets, we first compute the intersection-over-union (IoU) between detections and all predicted boxes of each target. Then we use the IoU matrix as the cost matrix of the Hungarian Algorithm for bipartite matching. If a target is not detected in 3 consecutive frames, we will delete the target from the scene.

### F. Future Trajectory Prediction

Trajectory prediction aims at predicting the future locations of objects. For an object, With the historical object positions obtained by the object tracking module

$$\mathbf{x}^{(t)} = (x_c^{(t)}, y_c^{(t)}) \tag{5}$$

$$\mathbf{X} = (\mathbf{x}^{(t-5)}, \mathbf{x}^{(t-4)}, ..., \mathbf{x}^{(t)}) \tag{6}$$

We feed the historical positions of all objects $\mathbf{X}_1, \mathbf{X}_1, ..., \mathbf{X}_n$ into a transformer $t$ [72], and predict the future positions of the objects. As shown in Fig 11, the transformer consists of multiple transformer encoder layers. Each transformer encoder layer contains a self-attention layer, a fully-connected (FC) layer, an Add + Norm layer, and a multi-layer perception (MLP) layer.

*a) Input encoding:* With the input historical trajectories of objects $\mathbf{X}_1, \mathbf{X}_2, ..., \mathbf{X}_n$, first we apply a positional mapping [15], [72] to obtain a higher dimensional vector and enable the prediction to easily approximate a higher frequency function. for each object historical trajectory $\mathbf{X}$, the positional mapping $\mathcal{T}$ is

$$\mathcal{T} : \mathbf{X} \mapsto (\mathbf{X}, \sin 2^0 \pi \mathbf{X}, \cos 2^L \pi \mathbf{X}, ..., \sin 2^0 \pi \mathbf{X}, \cos 2^L \pi \mathbf{X}) \tag{7}$$
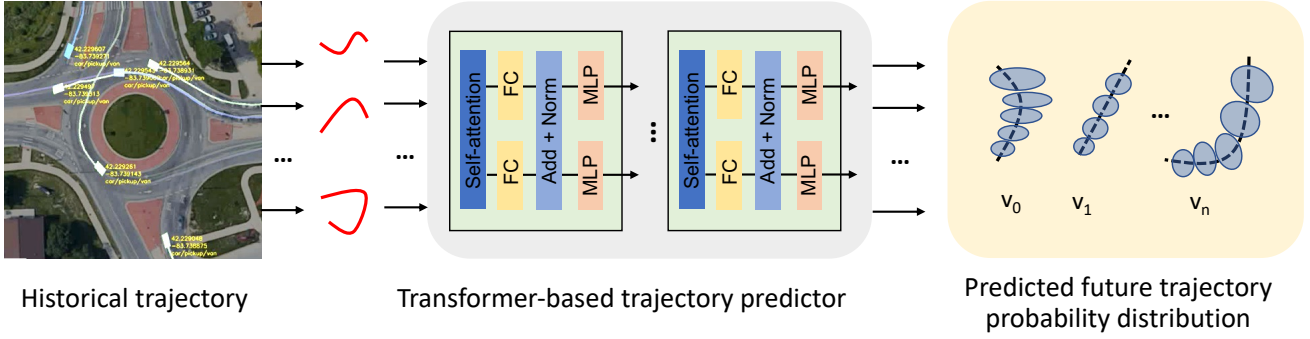
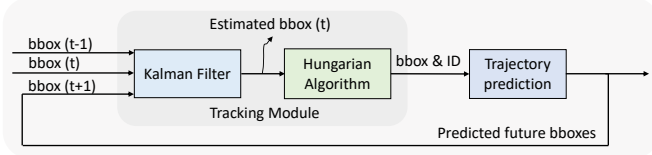Fig. 11: An illustration of the object prediction pipeline



Fig. 12: An illustration of the object tracking pipeline

$L$ is the positional mapping length parameter of $\mathcal{T}$. The output vector $\mathcal{T}\mathbf{X}$ is a vector with length $2L + 1$.

*b) Transformer-based trajectory predictor:* The trajectory predictor contains a linear layer to increase the dimension of the vectors to 256, a transformer with 4 encoder layers, and a predictor to predict trajectory mean and variance. In each encoder layer, there is a multi-head self-attention layer with the head number 8. The self-attention allows each object trajectory to interact with other object trajectories. This enables the predictor to predict complex driving behaviors like yielding to other vehicles. The MLP layer allows interactions within each object across different frames. The MLP layer contains a FC layer, a GeLU layer [73] and another FC layer. For the prediction head, two linear layers are used to predict future object location mean and variance accordingly.

*c) Output encoding:* The trajectory predictor predicts the object positions in the future 5 frames. The interval between two frames is 0.4 seconds. For each frame, we model the possible object position distribution as a Gaussian distribution $\mathcal{N}(\mu, \sigma^2)$. The mean $\mu$ and variance $\sigma^2$ are predicted.

*d) Loss function:* We use regression loss with uncertainty estimation as the training loss of the trajectory predictor. For each future frame and each object, we denote the object future position as $(x, y)$, the predicted object position mean as $(\mu_x, \mu_y)$, and predicted object position variance as $(\sigma_x, \sigma_y)$. The losses for the predictor training are

$$\ell(\mu) = (x - \mu_x)^2 + (y - \mu_y)^2 \tag{8}$$
$$\ell(\sigma_x^2) = ((x - \mu_x)^2 - \sigma_x^2)^2 \tag{9}$$
$$\ell(\sigma_y^2) = ((y - \mu_y)^2 - \sigma_y^2)^2 \tag{10}$$
$$\mathcal{L} = \ell(\mu) + \ell(\sigma_x^2) + \ell(\sigma_y^2) \tag{11}$$

## V. EXPERIMENTAL EVALUATIONS

### A. Experiment setup

**Training dataset.** The training dataset contains $4,000$ images, $1,000$ images for each camera (NE, NW, SE, SW). The images are collected in May and June of 2021. There are three object categories in the training dataset: cars, trucks/trailers/buses, and other road users.

**Training setup.** For object detection, we use the training pipeline of YOLOX [25]. We train the YOLOX-nano model for 150 epochs with a mini-batch size 8.The initial learning rate is $5e-5$ and decayed with a factor of 10 after 100 epochs. The weight decay is set to be $5e-4$ and Adam [74] optimizer is used. The model is trained on PyTorch [75] 1.9 platform with a NVIDIA RTX 3080 GPU. The input image size is $640 \times 640$. The data augmentation we used in training is the same as YOLOX's: random resize, horizontal flip, rotation, translation, shearing, and color distortion are applied.

**Evaluation setup.** We evaluate our infrastructure-based perception system at the same roundabout where we installed the system. The evaluation experiment proceeded on a Friday afternoon in October 2022 with a heavy traffic volume. A connected automated vehicle[1] – a Hybrid Lincoln MKZ 14, which is equipped with a high-precision RTK unit and an IMU is used for evaluation. The evaluation contains 8 trips at the roundabout: 4 trips contain a 270-degree turn using the inner lane, 4 trips contain a 90-degree turn using the outer lane. The routes of the trips are illustrated in Figure 13. The evaluation area is set to be the circular area around the roundabout within a 50 m radius. The sampling frequency of the RTK is 50 Hz, and the RTK latency is negligible. We use the vehicle trajectory produced by RTK as ground truth, and compare it with the trajectory detected from the perception system.

### B. Perception evaluation metrics

In the following, we describe the performance metrics of the MSight perception system.

**Detection metrics.** For object detection, use standard evaluation metrics including False Negative rate, False Positive rate, lateral position error, and longitudinal position error. Detections with lateral position errors smaller than 1.5 m
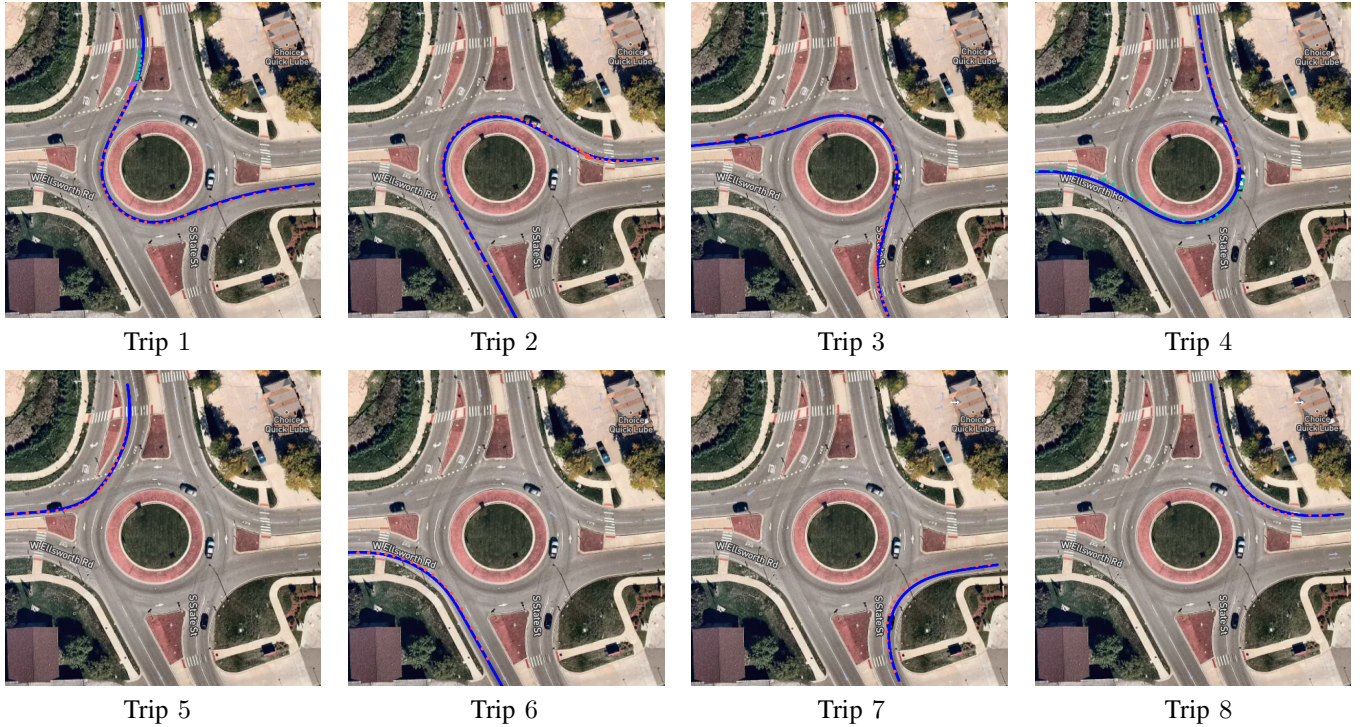
---

[1]https://mcity.umich.edu/

Fig. 13: Visualizations of the 8 trips in the field-test. The blue lines are the vehicle's ground-truth trajectories obtained by RTK. The colored dots are the detected vehicle locations obtained by MSight perception system with 4 cameras setup. Different colors mean different tracking IDs.
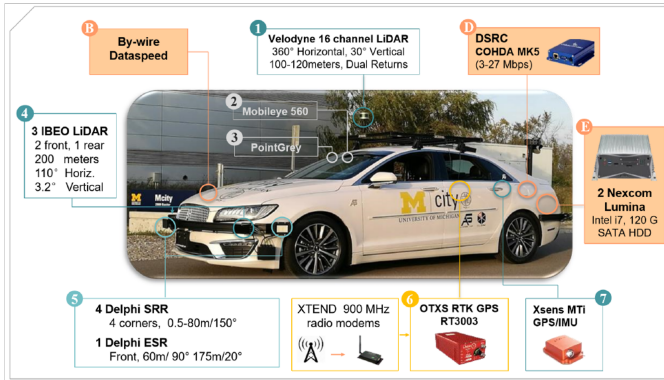


Fig. 14: An illustration of the vehicle for the field test. The RTK GPS, IMU, and V2X module are used for perception and system latency evaluations.



Fig. 15: An illustration of detection and tracking errors, including False Negative, False Positive, and ID Switch.

are regarded as True Positives, otherwise False Positives. The False Positive rate is defined as

$$\text{FP rate} = \frac{|\text{FP}|}{|\text{Dets}|} \qquad (12)$$

Here $|\text{Dets}|$ means the total number of detections. Given the detection frequency (in our case, 2.5 Hz), the expected timestamps and locations of the detections can be obtained. For each expected timestamp, if the detection at the timestamp is missing, or the detection lateral position error is larger than

1.5 m, a False Negative at this timestamp will be collected.

$$\text{FN rate} = \frac{|\text{FN}|}{|\text{gtDet}|} \qquad (13)$$

Here $|\text{gtDets}|$ means the expected number of ground-truth detections of the trip. For the lateral position error, we report the mean lateral position error of the True Positive detections.

**Tracking metrics.** For object tracking, we report the number of ID switches (#ID Switch), longest track preservation time (Longest Track (%)), and Multi-Object Tracking Accuracy (MOTA). The #ID Switch means how many times the tracking ID switches to another one for the same object. For example, in Figure 15, at first the object is tracked with ID 1, then ID switched to 2. In this case, the #ID Switch is 1. The longest

| Camera | NE | NW | SE | SW |
|---|---|---|---|---|
| Calibration error | 0.53 m | 0.47 m | 0.39 m | 0.46 m |

TABLE I: Calibration errors of four cameras.

| Date | Time | Error (w/o align) | Error (w/ align) |
|---|---|---|---|
| June 23th, 2022 | 10 am | 0.46 m | 0.46 m |
| June 23th, 2022 | 12 am | 0.57 m | 0.46 m |
| June 23th, 2022 | 2 pm | 0.62 m | 0.46 m |
| June 23th, 2022 | 4 pm | 0.63 m | 0.46 m |
| June 23th, 2022 | 6 pm | 0.65 m | 0.47 m |
| July 23th, 2022 | 10 am | 0.54 m | 0.46 m |
| August 23th, 2022 | 10 am | 2.97 m | 0.46 m |

TABLE II: Evaluations on the camera calibration and image alignment. The calibration is performed on June 23th at 10 am. This calibration has a larger error in other dates or times when without image alignment. With image alignment, the calibration errors are consistently small across different dates and times.

track preservation time means the number of frames of the longest track divided by the number of frames of the whole trip. In Figure 15, the Longest Track is 50%.

MOTA is one of the most representative measures of multi-object tracking. MOTA measures three types of errors mentioned above: False Positives, False Negatives, and #ID Switch. The MOTA is defined as

$$\text{MOTA} = 1 - \frac{|\text{FN}| + |\text{FP}| + |\text{IDSwitch}|}{|\text{gtDets}|} \quad (14)$$

For perfect detection and tracking, the MOTA score should be 1. The larger the MOTA score is, the fewer errors the model makes.

**Trajectory prediction metrics.** For trajectory prediction, we use standard Final Displacement Error (FDE) metrics. Suppose we want to predict the object location in next $K$ frames. The predicted object locations are $(\mathbf{x}_1, ..., \mathbf{x}_K)$, and the ground-truth object locations are $(\hat{\mathbf{x}}_1, ..., \hat{\mathbf{x}}_K)$. The FDE is defined as

$$\text{FDE}_K(m) = m(\mathbf{x}_K - \hat{\mathbf{x}}_K) \quad (15)$$

FDE measures the displacement between the predicted location and ground-truth location at the final frame. Here $m$ is an error measure. We use latitude position error and longitudinal error as the error measures. In our experiments, we set $K = 3$. With an interval of 0.4 seconds between frames, we predict the object locations in the future 1.2 seconds.

### C. Camera calibration and image alignment

**Camera calibration.** For each camera, we label about 20 landmarks both on the GoogleMap satellite image and the camera-captured image. Table I shows the average calibration error of the labeled landmarks. For all four cameras, the calibration errors are controlled within 0.6 m, which ensures the localization accuracy of our perception system.

**Image alignment.** As mentioned, roadside cameras will serve for a long time in their life cycle. The pose of the cameras might change slightly due to pole deformation. The pose changes will lead to inaccurate calibrations at different dates and times. To compensate for the issue, we develop an image alignment algorithm to align the images to a set of standard images. This will preserve consistent calibration quality across dates and times. As shown in Table II, we provide the calibration error at different dates and times. The calibration is done at 10 am, June 23th, 2022, and we re-evaluate the calibration quality at 12 pm, 2 pm, 4 pm, and 6 pm of the same day, as well as the same time but on July 23th, 2022, and August 23th, 2022. A change of date or time can lead to a larger calibration error. While, with image alignment, the calibration error is controlled to around 0.46 m, which is the original calibration error without change of date or time.
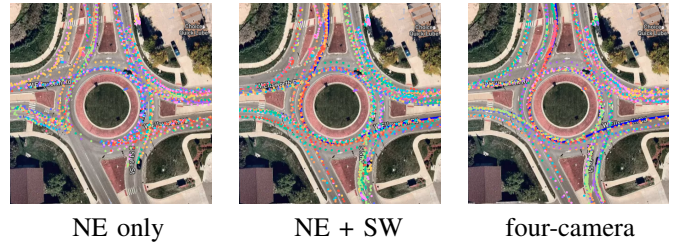


NE only      NE + SW      four-camera

Fig. 16: Visualizations of detection point distributions of different camera settings. NE only setting is not able to cover the whole roundabout. NE + SW and four-camera setting both have good coverage.

### D. Object detection and tracking results

Figure 13 and Table III show the perception results of our field test. In Figure 13, for all 8 trips, the localization accuracy of the MSight perception system is satisfactory. The colored dots (detected points) are overlapped with or very near to the blue lines (ground truth). For trip 1 and trip 4, there is one ID switch. We check the raw videos for these two trips, and find out that, the testing vehicle is occluded by a truck in the outer lane of the roundabout, and causes the miss detection and ID switch. We believe that this issue can be alleviated when a more advanced tracking algorithm is used. Further, the coverage of the MSight perception system is generally good. The blue lines are within the region-of-interest (within 50 m radius of the roundabout center). As shown in Figure 13, MSight has perfect horizontal coverage, while in the vertical direction, MSight has small blind spots at the top and bottom areas.

In Table III, we show the quantitative evaluation results of MSight perception system with different camera setups. There are three setups we used: (1) all four cameras at the northeast, northwest, southeast, southwest corner of the roundabout, (2) only two cameras in the diagonal direction: northeast and southwest corners, and (3) only one camera at the northeast corner.

For the four-camera setup, the MSight has the best detection quality with FN rate 11.83%, FP rate 4.51% and lateral error 0.63 m. However, in the four-camera setup, the MSight has a higher #ID Switch and shorter Longest Track. The

| Method | Camera Setup | Detection Evaluations | | | | Tracking Evaluations | | |
|--------|-------------|----------------|---------------|-----------------|-----------------|-------------|---------------------|-------|
| | | FN rate↓ (%) | FP rate↓ (%) | Lat. Error↓ (m) | Lon. Error↓ (m) | #ID Switch↓ | Longest Track↑ (%) | MOTA↑ |
| MSight | NE + NW + SE + SW | 11.83 | 4.51 | 0.63 | 0.90 | 0.38 | 90.25 | 0.82 |
| MSight | NE + SW | 22.26 | 9.69 | 0.63 | 1.03 | 0.00 | 100.00 | 0.66 |
| MSight | NE | 37.04 | 30.48 | 0.74 | 1.13 | 0.13 | 97.50 | 0.42 |

TABLE III: Evaluations of the MSight perception system. Three camera setups are tested: all four cameras, two cameras in diagonal directions, and only one camera. The camera setup with four cameras achieves overall the best performance. The reported results are the mean results of 8 trips.

| #Trip | FP rate | $FDE_{1.2s}$(Lat. error) | $FDE_{1.2s}$(Lon. error) |
|-------|---------|-----------------------|-----------------------|
| 1 | 20.26 | 0.79 m | 1.48 m |
| 2 | 2.99 | 0.74 m | 0.86 m |
| 3 | 4.08 | 0.62 m | 1.73 m |
| 4 | 15.79 | 0.60 m | 1.22 m |
| 5 | 36.36 | 0.75 m | 1.72 m |
| 6 | 19.05 | 0.32 m | 0.98 m |
| 7 | 35.00 | 0.84 m | 0.75 m |
| 8 | 4.55 | 0.82 m | 1.23 m |
| Overall | 16.01 | 0.69 m | 1.25 m |

TABLE IV: Evaluations of trajectory prediction performance of MSight system. The False Positive rate, latitude error and longitudinal error of the predicted vehicle positions are reported.

reason might be: more cameras cause more frequent switches between cameras, which might lead to more #ID Switches. MSight with four-camera achieves the best MOTA score 0.82. For the two-camera (NE+SW) setup, the detection quality is worse than the four-camera setup. The reason might be: with only two cameras, each camera needs to cover a larger area compared to the four-camera setup, and in the far-away areas, the localization is more inaccurate. For the one-camera (NE) setup, the detection quality is the worst. Both FN rate and FP rate exceed 30%. We find that only one camera is not able to cover the whole roundabout. As shown in Figure 16, in NE only figure, the detection density at the bottom-left area is clearly lower than other areas. It indicates that the bottom-left area is not well-covered. For NE+SW setup and four-camera setup, the coverage on the whole roundabout is good.

### E. Trajectory prediction results

Table IV shows the trajectory prediction results of the MSight system. We use the four-camera setup for this experiment to get the most accurate object detection results. The overall False Positive rate is 16.01, the $FDE_{1.2s}$(Lat. error) is 0.69 m, and $FDE_{1.2s}$(Lon. error) is 1.25 m. For the 8 trips, the trip #5 and trip #7 have a much larger error than other trips. We visualize the predicted trajectory of these two trips and find out that, our algorithm predicts a different path from the ground-truth path. In trip #5, the vehicle enters the roundabout in the outer lane (not the right-turn-only lane). Still, the vehicle performs a right turn, which is abnormal behavior. Our prediction algorithm predicts the vehicle to go straight instead, which is the most common driving behavior.

For trip #7, the prediction algorithm also falsely predict the future vehicle path.

### F. System latency

Latency is a very important aspect of cooperative driving as most cooperative driving systems are delay-sensitive. To analyze the system latency clearly, we break down the cooperative perception pipeline into **three** phases:

0) **Sensor data processing**: Sensors sense and generate the data. In this phase, the sensors sense the real-world environment and provide the raw sensor data (In this case, the fisheye images), we denote this phase **phase 0**.

1) **Perception algorithms execution**: In this phase, the perception algorithms introduced in Section IV are executed, and produce final perception results. This is the core perception phase of the overall cooperative perception pipeline, we denote this phase as **phase 1**.

2) **Perception results transmission**: In this phase, the perception results are transmitted to the vehicle from the roadside. During this phase, the perception results are first encoded into V2X messages; then, V2X radio forwards the encoded messages to the vehicle; finally, on the vehicle side, the onboard unit decodes the messages. We denote this phase as **phase 2**.

**Phase 0** is specific to sensors, which is not the focus of this paper. Consequently, in this paper, only latency in **phase 1** and **phase 2** is measured and presented. Figure 17 shows the breakdown of the perception pipeline and the measured latency.

The **phase 1** latency is crucial for determining how well the core perception algorithm performs. It is measured by running the algorithm in the production environment. Phase 1's average latency is **35 ms**, which proves the effectiveness of the core algorithms. As for **phase 2**, latency is measured by timestamping each message before encoding and comparing it with the time after decoding. **Phase 2** latency is also critical as this phase is the major additional component of cooperative perception in comparison to onboard perception. The average latency of phase 2 is **40 ms**. This indicates that the roadside cooperative perception system introduced in this paper adds only **40 ms** of latency over onboard perception (if all the other components execute equally fast). This proves that the system is highly viablefor cooperative driving in terms of latency.

### VI. CONCLUSION

In this paper, we introduce MSight, a comprehensive roadside cooperative perception system leveraging roadside cameras, designed explicitly for applications related to autonomous
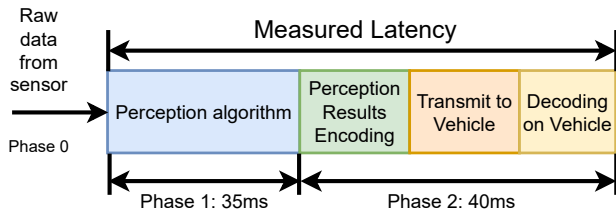
Fig. 17: Latency measurement result of the system

driving. MSight aims to augment onboard perception systems, addressing their limitations in complex situations where their performance is often suboptimal. Subsequently, the paper outlines field tests conducted in a production environment and discusses the derived results. The findings illustrate that the system can discern vehicles with lane precision and only imposes an additional 40 ms of latency to the onboard perception pipeline. These outcomes underscore the efficacy of the roadside perception system, revealing it as a viable solution for the development of cooperative driving systems in the future.

## REFERENCES

[1] M. Schwall, T. Daniel, T. Victor, F. Favaro, and H. Hohnhold, "Waymo public road safety performance data," *arXiv preprint arXiv:2011.00038*, 2020.

[2] J. B. Kenney, "Dedicated short-range communications (dsrc) standards in the united states," *Proceedings of the IEEE*, vol. 99, no. 7, pp. 1162–1182, 2011.

[3] S. Chen, J. Hu, Y. Shi, Y. Peng, J. Fang, R. Zhao, and L. Zhao, "Vehicle-to-everything (v2x) services supported by lte-based systems and 5g," *IEEE Communications Standards Magazine*, vol. 1, no. 2, pp. 70–76, 2017.

[4] S. Chen, J. Hu, Y. Shi, L. Zhao, and W. Li, "A vision of c-v2x: Technologies, field testing, and challenges with chinese development," *IEEE Internet of Things Journal*, vol. 7, no. 5, pp. 3872–3881, 2020.

[5] T. Kitazato, M. Tsukada, H. Ochiai, and H. Esaki, "Proxy cooperative awareness message: an infrastructure-assisted v2v messaging," in *2016 Ninth International Conference on Mobile Computing and Ubiquitous Networking (ICMU)*. IEEE, 2016, pp. 1–6.

[6] S. A. A. T. Committee *et al.*, "V2x sensor-sharing for cooperative & automated driving," *SAE J3224. Available online: https://www. sae. org/servlets/works/committeeHome. do*, 2019.

[7] S. R. E. Datondji, Y. Dupuis, P. Subirats, and P. Vasseur, "A survey of vision-based traffic monitoring of road intersections," *IEEE transactions on intelligent transportation systems*, vol. 17, no. 10, pp. 2681–2698, 2016.

[8] D. Meng, O. Sayer, R. Zhang, S. Shen, H. Li, and H. X. Liu, "Roco: A roundabout traffic conflict dataset," *arXiv preprint arXiv:2303.00563*, 2023.

[9] T. Furuya and C. J. Taylor, "Road intersection monitoring from video with large perspective deformation," Ph.D. dissertation, University of Pennsylvania, 2014.

[10] S. Messelodi, C. M. Modena, and M. Zanin, "A computer vision system for the detection and classification of vehicles at urban road intersections," *Pattern analysis and applications*, vol. 8, no. 1, pp. 17–31, 2005.

[11] R. Zhang, D. Meng, L. Bassett, S. Shen, Z. Zou, and H. X. Liu, "Robust roadside perception for autonomous driving: an annotation-free strategy with synthesized data," *arXiv preprint arXiv:2306.17302*, 2023.

[12] N. Saunier and T. Sayed, "A feature-based tracking algorithm for vehicles in intersections," in *The 3rd Canadian Conference on Computer and Robot Vision (CRV'06)*. IEEE, 2006, pp. 59–59.

[13] C. Li, A. Chiang, G. Dobler, Y. Wang, K. Xie, K. Ozbay, M. Ghandehari, J. Zhou, and D. Wang, "Robust vehicle tracking for urban traffic videos at intersections," in *2016 13th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS)*. IEEE, 2016, pp. 207–213.

[14] F. Faisal, S. K. Das, A. H. Siddique, M. Hasan, S. Sabrin, C. A. Hossain, and Z. Tong, "Automated traffic detection system based on image processing," *Journal of Computer Science and Technology Studies*, vol. 2, no. 1, pp. 18–25, 2020.

[15] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," in *European conference on computer vision*. Springer, 2020, pp. 213–229.

[16] D. Meng, X. Chen, Z. Fan, G. Zeng, H. Li, Y. Yuan, L. Sun, and J. Wang, "Conditional detr for fast training convergence," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 3651–3660.

[17] X. Chen, F. Wei, G. Zeng, and J. Wang, "Conditional detr v2: Efficient detection transformer with box queries," *arXiv preprint arXiv:2207.08914*, 2022.

[18] X. Zhu, W. Su, L. Lu, B. Li, X. Wang, and J. Dai, "Deformable detr: Deformable transformers for end-to-end object detection," *arXiv preprint arXiv:2010.04159*, 2020.

[19] S. Liu, F. Li, H. Zhang, X. Yang, X. Qi, H. Su, J. Zhu, and L. Zhang, "Dab-detr: Dynamic anchor boxes are better queries for detr," *arXiv preprint arXiv:2201.12329*, 2022.

[20] Q. Chen, X. Chen, G. Zeng, and J. Wang, "Group detr: Fast training convergence with decoupled one-to-many label assignment," *arXiv preprint arXiv:2207.13085*, 2022.

[21] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 779–788.

[22] J. Redmon and A. Farhadi, "Yolo9000: better, faster, stronger," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 7263–7271.

[23] ——, "Yolov3: An incremental improvement," *arXiv preprint arXiv:1804.02767*, 2018.

[24] A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao, "Yolov4: Optimal speed and accuracy of object detection," *arXiv preprint arXiv:2004.10934*, 2020.

[25] Z. Ge, S. Liu, F. Wang, Z. Li, and J. Sun, "Yolox: Exceeding yolo series in 2021," *arXiv preprint arXiv:2107.08430*, 2021.

[26] D. Barina, "Gabor wavelets in image processing," *arXiv preprint arXiv:1602.03308*, 2016.

[27] J. Hariyono, V.-D. Hoang, and K.-H. Jo, "Moving object localization using optical flow for pedestrian detection from a moving vehicle," *The Scientific World Journal*, vol. 2014, 2014.

[28] D.-S. Kim and J. Kwon, "Moving object detection on a vehicle mounted back-up camera," *Sensors*, vol. 16, no. 1, p. 23, 2015.

[29] R. Yao, G. Lin, S. Xia, J. Zhao, and Y. Zhou, "Video object segmentation and tracking: A survey," *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 11, no. 4, pp. 1–47, 2020.

[30] H. A. Patel and D. G. Thakore, "Moving object tracking using kalman filter," *International Journal of Computer Science and Mobile Computing*, vol. 2, no. 4, pp. 326–332, 2013.

[31] L. E. Taylor, M. Mirdanies, and R. P. Saputra, "Optimized object tracking technique using kalman filter," *arXiv preprint arXiv:2103.05467*, 2021.

[32] S.-H. Bae and K.-J. Yoon, "Robust online multi-object tracking based on tracklet confidence and online discriminative appearance learning," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 1218–1225.

[33] J. Berclaz, F. Fleuret, E. Turetken, and P. Fua, "Multiple object tracking using k-shortest paths optimization," *IEEE transactions on pattern analysis and machine intelligence*, vol. 33, no. 9, pp. 1806–1819, 2011.

[34] P. Bergmann, T. Meinhardt, and L. Leal-Taixe, "Tracking without bells and whistles," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 941–951.

[35] C. Dicle, O. I. Camps, and M. Sznaier, "The way they move: Tracking multiple targets with similar appearance," in *Proceedings of the IEEE international conference on computer vision*, 2013, pp. 2304–2311.

[36] A. Milan, S. Roth, and K. Schindler, "Continuous energy minimization for multitarget tracking," *IEEE transactions on pattern analysis and machine intelligence*, vol. 36, no. 1, pp. 58–72, 2013.

[37] A. Bewley, Z. Ge, L. Ott, F. Ramos, and B. Upcroft, "Simple online and realtime tracking," in *2016 IEEE international conference on image processing (ICIP)*. IEEE, 2016, pp. 3464–3468.

[38] A. Aboah, "A vision-based system for traffic anomaly detection using deep learning and decision trees," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 4207–4212.

[39] L. Huang and W. Huang, "Rd-yolo: an effective and efficient object detector for roadside perception system," *Sensors*, vol. 22, no. 21, p. 8097, 2022.

[40] W. Wang, T. Gee, J. Price, and H. Qi, "Real time multi-vehicle tracking and counting at intersections from a fisheye camera," in *2015 IEEE Winter Conference on Applications of Computer Vision*. IEEE, 2015, pp. 17–24.

[41] A. Gaszczak, T. P. Breckon, and J. Han, "Real-time people and vehicle detection from uav imagery," in *Intelligent Robots and Computer Vision XXVIII: Algorithms and Techniques*, vol. 7878. International Society for Optics and Photonics, 2011, p. 78780B.

[42] Y. Iwasaki, M. Misumi, and T. Nakamiya, "Robust vehicle detection under various environmental conditions using an infrared thermal camera and its application to road traffic flow monitoring," *Sensors*, vol. 13, no. 6, pp. 7756–7773, 2013.

[43] P. Sun, C. Sun, R. Wang, and X. Zhao, "Object detection based on roadside lidar for cooperative driving automation: A review," *Sensors*, vol. 22, no. 23, p. 9316, 2022.

[44] Z. Zhang, J. Zheng, X. Wang, and X. Fan, "Background filtering and vehicle detection with roadside lidar based on point association," in *2018 37th Chinese Control Conference (CCC)*. IEEE, 2018, pp. 7938–7943.

[45] J. Zhao, H. Xu, H. Liu, J. Wu, Y. Zheng, and D. Wu, "Detection and tracking of pedestrians and vehicles using roadside lidar sensors," *Transportation research part C: emerging technologies*, vol. 100, pp. 68–87, 2019.

[46] Z. Zhang, J. Zheng, H. Xu, and X. Wang, "Vehicle detection and tracking in complex traffic circumstances with roadside lidar," *Transportation research record*, vol. 2673, no. 9, pp. 62–71, 2019.

[47] S. Zhou, H. Xu, G. Zhang, T. Ma, and Y. Yang, "Leveraging deep convolutional neural networks pre-trained on autonomous driving data for vehicle detection from roadside lidar data," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 11, pp. 22 367–22 377, 2022.

[48] Z. Bai, G. Wu, M. J. Barth, Y. Liu, E. A. Sisbot, and K. Oguchi, "Pillargrid: Deep learning-based cooperative perception for 3d object detection from onboard-roadside lidar," in *2022 IEEE 25th International Conference on Intelligent Transportation Systems (ITSC)*. IEEE, 2022, pp. 1743–1749.

[49] J. Wu, H. Xu, Y. Tian, R. Pi, and R. Yue, "Vehicle detection under adverse weather from roadside lidar data," *Sensors*, vol. 20, no. 12, p. 3433, 2020.

[50] J. Wu, H. Xu, J. Zheng, and J. Zhao, "Automatic vehicle detection with roadside lidar data under rainy and snowy conditions," *IEEE Intelligent Transportation Systems Magazine*, vol. 13, no. 1, pp. 197–209, 2020.

[51] M. Shan, K. Narula, Y. F. Wong, S. Worrall, M. Khan, P. Alexander, and E. Nebot, "Demonstrations of cooperative perception: safety and robustness in connected and automated vehicle operations," *Sensors*, vol. 21, no. 1, p. 200, 2021.

[52] R. Zhang, Z. Zou, S. Shen, and H. X. Liu, "Design, implementation, and evaluation of a roadside cooperative perception system," *Transportation Research Record*, p. 03611981221092402, 2022.

[53] Z. Zou, R. Zhang, S. Shen, G. Pandey, P. Chakravarty, A. Parchami, and H. X. Liu, "Real-time full-stack traffic scene perception for autonomous driving with roadside cameras," in *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022, pp. 890–896.

[54] Y. Du, B. Qin, C. Zhao, Y. Zhu, J. Cao, and Y. Ji, "A novel spatio-temporal synchronization method of roadside asynchronous mmw radar-camera for sensor fusion," *IEEE Transactions on Intelligent Transportation Systems*, 2021.

[55] J. Gao, C. Sun, H. Zhao, Y. Shen, D. Anguelov, C. Li, and C. Schmid, "Vectornet: Encoding hd maps and agent dynamics from vectorized representation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 11 525–11 533.

[56] H. Cui, V. Radosavljevic, F.-C. Chou, T.-H. Lin, T. Nguyen, T.-K. Huang, J. Schneider, and N. Djuric, "Multimodal trajectory predictions for autonomous driving using deep convolutional networks," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 2090–2096.

[57] M. Liang, B. Yang, R. Hu, Y. Chen, R. Liao, S. Feng, and R. Urtasun, "Learning lane graph representations for motion forecasting," in *European Conference on Computer Vision*. Springer, 2020, pp. 541–556.

[58] M. Wang, X. Zhu, C. Yu, W. Li, Y. Ma, R. Jin, X. Ren, D. Ren, M. Wang, and W. Yang, "Ganet: Goal area network for motion forecasting," *arXiv preprint arXiv:2209.09723*, 2022.

[59] W. Luo, B. Yang, and R. Urtasun, "Fast and furious: Real time end-to-end 3d detection, tracking and motion forecasting with a single convolutional net," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2018, pp. 3569–3577.

[60] D. Meng, C. Yu, J. Deng, D. Qian, H. Li, and D. Ren, "Hybrid motion representation learning for prediction from raw sensor data," *IEEE Transactions on Multimedia*, pp. 1–12, 2023.

[61] A. Rauch, F. Klanner, and K. Dietmayer, "Analysis of v2x communication parameters for the development of a fusion architecture for cooperative perception systems," in *2011 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2011, pp. 685–690.

[62] A. Rauch, F. Klanner, R. Rasshofer, and K. Dietmayer, "Car2x-based perception in a high-level fusion architecture for cooperative perception systems," in *2012 IEEE Intelligent Vehicles Symposium*. IEEE, 2012, pp. 270–275.

[63] M. Tsukada, M. Kitazawa, T. Oi, H. Ochiai, and H. Esaki, "Cooperative awareness using roadside unit networks in mixed traffic," in *2019 IEEE Vehicular Networking Conference (VNC)*. IEEE, 2019, pp. 1–8.

[64] M. Tsukada, T. Oi, A. Ito, M. Hirata, and H. Esaki, "Autoc2x: Open-source software to realize v2x cooperative perception among autonomous vehicles," in *2020 IEEE 92nd Vehicular Technology Conference (VTC2020-Fall)*. IEEE, 2020, pp. 1–6.

[65] M. Tsukada, T. Oi, M. Kitazawa, and H. Esaki, "Networked roadside perception units for autonomous driving," *Sensors*, vol. 20, no. 18, p. 5320, 2020.

[66] S. Yang, H. H. Yin, R. W. Yeung, X. Xiong, Y. Huang, L. Ma, M. Li, and C. Tang, "On scalable network communication for infrastructure-vehicle collaborative autonomous driving," *IEEE Open Journal of Vehicular Technology*, 2022.

[67] J. Kannala and S. S. Brandt, "A generic camera model and calibration method for conventional, wide-angle, and fish-eye lenses," *IEEE transactions on pattern analysis and machine intelligence*, vol. 28, no. 8, pp. 1335–1340, 2006.

[68] S. Shah and J. Aggarwal, "Intrinsic parameter calibration procedure for a (high-distortion) fish-eye lens camera with distortion model and accuracy estimation," *Pattern Recognition*, vol. 29, no. 11, pp. 1775–1788, 1996.

[69] G. D. Evangelidis and E. Z. Psarakis, "Parametric image alignment using enhanced correlation coefficient maximization," *IEEE transactions on pattern analysis and machine intelligence*, vol. 30, no. 10, pp. 1858–1865, 2008.

[70] G. Welch, G. Bishop *et al.*, "An introduction to the kalman filter," 1995.

[71] H. W. Kuhn, "The hungarian method for the assignment problem," *Naval research logistics quarterly*, vol. 2, no. 1-2, pp. 83–97, 1955.

[72] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.

[73] D. Hendrycks and K. Gimpel, "Gaussian error linear units (gelus)," *arXiv preprint arXiv:1606.08415*, 2016.

[74] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[75] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga *et al.*, "Pytorch: An imperative style, high-performance deep learning library," *Advances in neural information processing systems*, vol. 32, 2019.
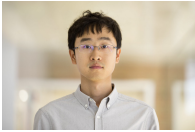
## VII. Biography

**Rusheng Zhang** received the B.E. degree in micro electrical mechanical system and second B.E. degree in Applied Mathematics from Tsinghua University, Beijing, in 2013. He received an M.S. and phD degree in electrical and computer engineering from Carnegie Mellon University, in 2015, 2019 respectively. His research areas include artificial intelligence, cooperative driving, cloud computing and vehicular networks.

**Depu Meng** (Member, IEEE) is a Post Doctoral Research Fellow at the Department of Civil and Environmental Engineering, University of Michigan. He received his B. E. degree from the Department of Electrical Engineering and Information Science at the University of Science and Technology of China in 2018. He received his Ph. D. degree from the Department of Automation at the University of Science and Technology of China. His research interests include computer vision and autonomous driving systems.

**Shengyin (Sean) Shen** works as a Research Engineer in the Engineering Systems Group at the University of Michigan Transportation Research Institute (UMTRI). Sean holds an MS degree in Civil and Environmental Engineering from the University of Michigan, Ann Arbor, and an MS degree in Electrical Engineering from the University of Bristol, UK. He also earned a BS degree from Beijing University of Posts and Telecommunications, China. Sean's research interests are primarily focused on cooperative driving automation and related applications that use roadside perception, edge-cloud computing, and V2X communications to accelerate the deployment of automated vehicles. He has extensive experience in implementation of large-scale deployments, such as the Safety Pilot Model Deployment (SPMD), Ann Arbor Connected Vehicle Testing Environment (AACVTE), and Smart Intersection Project. Moreover, he has been involved in many research projects funded by public agencies such as USDOT, USDOE, and companies such as Crash Avoidance Metric Partnership (CAMP), Ford Motor Company, and GM Company, among others.

**Zhengxia Zou** Zhengxia Zou received his B.S. and Ph.D. degrees from the Image Processing Center, School of Astronautics, Beihang University, Beijing, China, in 2013 and 2018, respectively. He is currently a Professor at the School of Astronautics, Beihang University. From 2018 to 2021, he worked at the University of Michigan, Ann Arbor as a Post-Doctoral Research Fellow. His research interests include computer vision and related problems in autonomous driving and remote sensing. He has published more than 40 peer-reviewed papers in top-tier journals and conferences, including Nature, Nature Communications, PROCEEDINGS OF THE IEEE, IEEE TRANSACTIONS ON IMAGE PROCESSING, IEEE TRANSACTIONS ON GEOSCIENCE AND REMOTE SENSING, and IEEE/CVF Computer Vision and Pattern Recognition. Dr. Zou was selected as "World's Top 2% Scientists" by Stanford University in 2022.

**Houqiang Li** (Fellow, IEEE) is a Professor with the Department of Electronic Engineering and Information Science at the University of Science and Technology of China. His research interests include multimedia search, image/video analysis, video coding and communication. He has authored and co-authored over 200 papers in journals and conferences. He is the winner of National Science Funds (NSFC) for Distinguished Young Scientists, the Distinguished Professor of Changjiang Scholars Program of China, and the Leading Scientist of Ten Thousand Talent Program of China. He served as an Associate Editor of the IEEE Transactions on Circuits and Systems for Video Technology from 2010 to 2013. He served as the TPC Co-Chair of VCIP 2010, and he served as the General Co-Chair of ICME 2021. He is the recipient of National Technological Invention Award of China (second class) in 2019 and the recipient of National Natural Science Award of China (second class) in 2015. He was the recipient of the Best Paper Award for VCIP 2012, the recipient of the Best Paper Award for ICIMCS 2012, and the recipient of the Best Paper Award for ACM MUM in 2011.

Houqiang received the B.S., M. Eng. and Ph.D degrees in electronic engineering from the University of Science and Technology of Chinae, Hefei, China in 1992, 1997 and 2000, respectively. He was elected as a Fellow of IEEE (2021).

**Henry X. Liu** (Member, IEEE) received the bachelor's degree in automotive engineering from Tsinghua University, China, in 1993, and the PhD. degree in civil and environment engineering from the University of Wisconsin-Madison in 2000. He is currently a professor in the Department of Civil and Environmental Engineering and the Director of Mcity at the University of Michigan, Ann Arbor. He is also a Research Professor at the University of Michigan Transportation Research Institute and the Director for the Center for Connected and Automated Transportation (USDOT Region 5 University Transportation Center). From August 2017 to August 2019, Prof. Liu served as DiDi Fellow and Chief Scientist on Smart Transportation for DiDi Global, Inc., one of the leading mobility service providers in the world. Prof. Liu conducts interdisciplinary research at the interface of transportation engineering, automotive engineering, and artificial intelligence. Specifically, his scholarly interests concern traffic flow monitoring, modeling, and control, as well as testing and evaluation of connected and automated vehicles. Prof. Liu is the managing editor of Journal of Intelligent Transportation Systems.