

# Computer Vision Aided mmWave Beam Alignment in V2X Communications

Weihua Xu, Feifei Gao, Xiaoming Tao, Jianhua Zhang, and Ahmed Alkhateeb

## Abstract

Visual information, captured for example by cameras, can effectively reflect the sizes and locations of the environmental scattering objects, and thereby can be used to infer communications parameters like propagation directions, receiver powers, as well as the blockage status. In this paper, we propose a novel beam alignment framework that leverages images taken by cameras installed at the mobile user. Specifically, we utilize 3D object detection techniques to extract the size and location information of the dynamic vehicles around the mobile user, and design a deep neural network (DNN) to infer the optimal beam pair for transceivers without any pilot signal overhead. Moreover, to avoid performing beam alignment too frequently or too slowly, a beam coherence time (BCT) prediction method is developed based on the vision information. This can effectively improve the transmission rate compared with the beam alignment approach with the fixed BCT. Simulation results show that the proposed vision based beam alignment methods outperform the existing LIDAR and vision based solutions, and demand for much lower hardware cost and communication overhead.

## Index Terms

Deep learning, beam alignment, beam coherence time, computer vision, V2X communication

W. Xu and F. Gao are with Institute for Artificial Intelligence, Tsinghua University (THUAI), Beijing National Research Center for Information Science and Technology (BNRist), Department of Automation, Tsinghua University, Beijing, P.R. China, 100084 (email: xwh19@mails.tsinghua.edu.cn, feifeigao@ieee.org).

X. Tao is with the Department of Electronic Engineering, Tsinghua University, Beijing, P.R. China, 100084 (email: taoxm@tsinghua.edu.cn).

J. Zhang is with the State Key Laboratory of Networking and Switching Technology, Beijing University of Posts and Telecommunications, Beijing 100876, China (e-mail: jhzhzhang@bupt.edu.cn).

A. Alkhateeb is with the Department of Electrical and Computer Engineering, University of Texas at Austin, Austin, TX 78712-1687 USA (e-mail: alkhateeb@asu.edu).

## I. INTRODUCTION

Beamforming has been deemed as the critical technique to overcome the signal attenuation of high frequency band, especially for the millimeter wave or even Terahertz communication [1]. Traditional beamforming strategies include sweeping the beam with a certain codebook [2] to maximize the receive signal-to-noise ratio (SNR) or calculating the beamforming matrix directly from the estimated channel matrix [3]. However, it is generally known that producing high gain beam with large antenna array leads to the huge time and spectrum overhead. Recently, integrated sensing and communication (ISAC) has drawn great attention for its capability to assist beamforming [4]-[18]. The principle behind is that the sensor data, i.e., the point cloud or the RGB/depth images from GPS, Radar, LIDAR, and camera equipped at many intelligent terminals, are the valid out-of-band information that can indicate the spatial characteristics of the communications environment.

In [6], the authors proposed to use the locations of the mobile station (MS) and the surrounding vehicles to realize power estimation of all the beam pairs under the V2X scenario. Although the high throughput ratio can be achieved by [6], the requirement for all surrounding vehicles to transmit the instantaneous locations to the MS will introduce the huge cost and delay. Thus, the *proactive* perception capability of the BS is further explored to reduce the transmission overhead of surrounding vehicles as well as to obtain more accurate and abundant environmental information. With the aid of the Radar mounted at BS, [7] proposes a new hybrid beamforming scheme, and [8] presents a beam tracking method in multi-user situation. The authors in [9] propose a Radar-aided beam alignment method, and verify it by the real-world dataset. Nevertheless, the detection accuracy of Radar is inferior to the LIDAR and camera, as the Radar signal is mainly suitable to detect the directional information of transceivers but cannot accurately reflect the precise distribution of the surrounding scatters. In [10]-[11], the Base Station (BS) and MS use LIDAR to scan point clouds for accelerating the beam alignment. In [12], the authors proposed a federated learning based beam alignment method to use the LIDAR data from multiple distributed MSs and train a DNN cooperatively for realizing better beam alignment accuracy. However, LIDAR is expensive and is not universally implemented compared to the ordinary cameras.

The authors in [13] then proposed to equip the BS with multiple cameras and take the RGB images of the user and surrounding vehicles. The beam sequence and the image coordinates

of environmental objects are utilized as the input feature to predict blockage and BS handover from the recurrent neural network (RNN). In [14]-[15], the previous images taken at BS and the beam sequences are utilized for beam tracking. In [16], a multi-modal beam alignment method is proposed by the conjunctive use of the images taken at BS and the location coordinates of MS.

Nevertheless, if BS is obliged to take images, then additional characteristic information of the MS is required by BS to identify the MS from all detected objects inside one image, i.e., to match the MS's radio signal with the corresponding visual information [17]. Therefore, in [13]-[16], the BS needs to rely on the previous beam sequences or the location fed back by MS to achieve the blockage/handover prediction or beam tracking/selection, which leads to inevitable communication overhead. When the MS is blocked by the surrounding objects in the camera view of BS, the visual information of the MS will be lost to cause the degradation of the beamforming performance. Moreover, taking image at BS may violate the privacy of the customers and may not be accepted in real implementation.

Thus, an alternative way is to utilize the vision capability of the intelligent MS, such as autonomous car and unmanned aerial vehicle, to avoid the overhead for MS identification and obtain better beamforming performance. In fact, the visual resources at MS are more convenient to be obtained, especially with the gradual popularization of automatic driving, and UAV reconnaissance, since the visual data are also widely used for navigation and obstacle avoidance. In [18], the authors use the scene images taken at MS to infer the channel covariance matrix. However, the covariance matrix is only the statistical characteristic of channel, whereas the instantaneous beam alignment relying on the vision of user terminal has not been tackled yet, to the best of the authors' knowledge.

In this paper, we propose to utilize the camera images taken at MS for the beam alignment under dynamic environment. Specifically, the main contributions of this paper contains following three different aspects:

- 1) Vision based beam alignment when the MS location is available (VBALA), in which the 3D detection technique is utilized to obtain the 3D spatial distribution information of MS surrounding dynamic objects, and then a DNN is designed to predict the optimal beam pair from the MS location and vehicle distribution information. Vision based beam alignment when the MS location is unavailable or not accurate enough due to the strong noise corruption (VBALU), in which the optimal beam is predicted only through the images.

TABLE I  
THE RELATED RESEARCH WORKS FOR COMPARISON

Paper	Sensor	Sensor Data Type	Task	Limitations compared with the proposed methods
[6]	GPS at MS and other vehicles	Location coordinates	Beam power estimation	Huge cost and delay for transmission of instantaneous locations
[7],[9]	Radar at BS	Radar echo signals	Beam alignment	Inferior detection accuracy than the camera
[8]			Beam tracking	
[10]	LIDAR at BS or MS	Point Clouds	Beam alignment	More expensive and less universal than the ordinary camera
[11]-[12]	LIDAR at MS			
[13]	Camera at BS	RGB Images	Blockage prediction and BS handover prediction	Requirement for MS identification and privacy concerns
[14]-[15]			Beam tracking	
[16]			Beam alignment	
[18]	Camera at MS		Channel covariance estimation	Only statistical characteristic estimation of channel
Ours			Beam alignment and BCT prediction	

- 2) Moreover, we present a vision based method to predict the beam coherent time (BCT) (VPBCT), i.e., the duration that the optimal beam pair remains unchanged for codebook-based beamforming.
- 3) By the image and channel data generated from the autonomous driving and ray tracing simulation software, the proposed VBALA and VBALU are verified to achieve better beam alignment performance than the Lidar and BS's vision based methods, which demonstrates the advantage of utilizing the images taken at MS. The proposed VPBCT can outperform the fixed BCT based beamforming approach. The simulated dataset is made available to the public [20].

The above related research works are listed in the TABLE I for clear comparison with the proposed methods.

This paper is organized as follows. Section II introduces the signal model as well as the beam alignment criteria of the communication system. Section III proposes the vision based beam alignment methods with/without MS location, while Section IV presents the vision based BCT prediction method. Section V provides the performance metric, the simulation setup, and numerical results together with detailed discussions. Finally, Section VI draws the conclusions.

TABLE II  
CRITICAL NOTATIONS IN THE PAPER

Notation	Description	Notation	Description
$N_B$ and $N_U$	Antenna numbers of BS and MS	$\mathbf{F}$	Vehicle distribution feature for VBALA
$K$	Subcarrier number	$\mathbf{I}$	Scene image feature for VBALU
$\mathbf{w}_B$ and $\mathbf{w}_U$	Transmit and receive beamforming vector	$\mathbf{D}$	Sequence of scene image features for VPBCT
$N_B^{CB}$ and $N_U^{CB}$	Codebook sizes of transmit and receive beam	$T_d$	Shooting interval of the camera
$C$	Number of cameras equipped at MS	$T_b$	Time for beam alignment during one BCT
$\theta_M^i$	Azimuth of the $i$ th camera	$M$	Ratio between BCT and shooting interval
$L_G$ and $W_G$	Length and width of the grid for VBALA	$\mathcal{Q}_T, \mathcal{Q}_V$ and $\mathcal{Q}_E$	Index set of image set sequences for constructing training, validation and test set

Notation:  $\mathbf{A}$  is a matrix;  $\mathcal{A}$  is a set;  $\mathbf{a}$  is a vector;  $a$  is a scalar;  $\mathbf{A}_{[i,j]}$  is the element of the  $i$ th row and the  $j$ th column in  $\mathbf{A}$ ;  $\mathbf{A}_{[i,:]}$  and  $\mathbf{A}_{[:,j]}$  are the  $i$ th row and the  $j$ th column of  $\mathbf{A}$  respectively;  $\mathcal{N}(\mathbf{m}_g, \mathbf{R}_g)/\mathcal{CN}(\mathbf{m}_g, \mathbf{R}_g)$  is the real/complex Gaussian random distribution with mean  $\mathbf{m}_g$  and covariance  $\mathbf{R}_g$ ;  $\text{Card}(\mathcal{A})$  is the cardinality of the set  $\mathcal{A}$ ;  $\mathbb{E}\{\cdot\}$  is the expectation operator. For the convenience of expression and reference, the critical notations adopted in the paper are summarized in TABLE II.

## II. SIGNAL MODEL

Let us consider a downlink orthogonal frequency-division multiplexing (OFDM) mmWave communication system with a single BS and a single MS. BS is equipped with a uniform linear array (ULA) of  $N_B$  antennas, and the MS is equipped with a ULA of  $N_U$  antennas. Both the BS and MS are assumed to have only a radio frequency chain. The downlink signal received at the user for the  $k$ th subcarrier can be expressed as

$$y_k = \mathbf{w}_U^H \mathbf{H}_k \mathbf{w}_B s_k + n_k, \quad (1)$$

where  $s_k \in \mathbb{C}$  is the transmit signal,  $\mathbf{H}_k \in \mathbb{C}^{N_U \times N_B}$  is the downlink channel matrix at the  $k$ th subcarrier,  $\mathbf{w}_B \in \mathbb{C}^{N_B \times 1}$  is the transmit beamforming vector,  $\mathbf{w}_U \in \mathbb{C}^{N_U \times 1}$  is the receive beamforming vector, and  $n_k \in \mathcal{CN}(0, \sigma^2)$  is the Gaussian noise at the  $k$ th subcarrier. The transmit signal  $s_k$  satisfies  $\mathbb{E}\{|s_k|^2\} = P_k$ .

From the widely used geometric channel model [19], the channel matrix  $\mathbf{H}_k$  can be expressed as

$$\mathbf{H}_k = \sum_{n=0}^{N-1} \sum_{l=1}^L \alpha_l e^{-j \frac{2\pi k}{K} n} d(nT_s - \tau_l) \mathbf{a}_r(\phi_l^r) \mathbf{a}_t^H(\phi_l^t), \quad (2)$$

where  $\alpha_l$  is the complex gain of the  $l$ th path,  $\tau_l$  is the time delay of the  $l$ th path,  $\phi_l^r$  and  $\phi_l^t$  are the  $l$ th path's azimuth angle of arrival and departure respectively,  $T_s$  is the sampling interval,  $K$  is the number of subcarriers,  $N$  is the length of cyclic prefix, and  $d(\cdot)$  denotes the pulse shaping filter. Moreover,  $\mathbf{a}_r(\phi) \in \mathbb{C}^{N_U \times 1}$  and  $\mathbf{a}_t(\phi) \in \mathbb{C}^{N_B \times 1}$  are the complex steering vector of the receive and transmit ULA array respectively.

When the antenna spacing is set as the half carrier wavelength, the mathematical expression of steering vector  $\mathbf{a}_r(\phi)$  and  $\mathbf{a}_t(\phi)$  is

$$\begin{aligned} \mathbf{a}_r(\phi) &= \frac{1}{\sqrt{N_U}} [1, e^{j\pi \sin(\phi)}, \dots, e^{j(N_U-1)\pi \sin(\phi)}]^T, \\ \mathbf{a}_t(\phi) &= \frac{1}{\sqrt{N_B}} [1, e^{j\pi \sin(\phi)}, \dots, e^{j(N_B-1)\pi \sin(\phi)}]^T. \end{aligned} \quad (3)$$

Assume the system utilizes the beamforming codebook to perform analog beamforming. The optimal beam pair  $(\mathbf{w}_B^{\text{opt}}, \mathbf{w}_U^{\text{opt}})$  should be selected from the transmit beam codebook  $\mathcal{W}_B = \{\mathbf{w}_{B,1}, \mathbf{w}_{B,2}, \dots, \mathbf{w}_{B,N_B^{\text{CB}}}\}$  and the receive beam codebook  $\mathcal{W}_U = \{\mathbf{w}_{U,1}, \mathbf{w}_{U,2}, \dots, \mathbf{w}_{U,N_U^{\text{CB}}}\}$  by maximizing the data rate, i.e.,

$$(\mathbf{w}_B^{\text{opt}}, \mathbf{w}_U^{\text{opt}}) = \arg \max_{(\mathbf{w}_B, \mathbf{w}_U)} \frac{1}{K} \sum_{k=1}^K \log_2 \left( 1 + \frac{P_k}{\sigma^2} |\mathbf{w}_U^H \mathbf{H}_k \mathbf{w}_B|^2 \right). \quad (4)$$

### III. VISION BASED BEAM ALIGNMENT

The concerned V2X scenario for vision based beam alignment is shown in Fig. 1. The vehicles equipped with many sensors for driver assistance are gradually popularized with the rapid development of autonomous driving technology [21]-[22]. We consider the MS is the target vehicle with  $C$  auxiliary cameras and runs along the traffic lane, i.e., we assume the MS movement direction is parallel to the traffic lane's right direction. Moreover, there are many other random dynamic vehicles traveling on the same or the nearby lane and acting as the possible blockage objects.

A road-side unit (RSU) acts as the BS and communicates with the MS. The RSUs are expected to be widely deployed in the traffic environments to realize the highly efficient vehicle to infrastructure (V2I) communication for the vehicle safety [23], the vehicle network [24] and



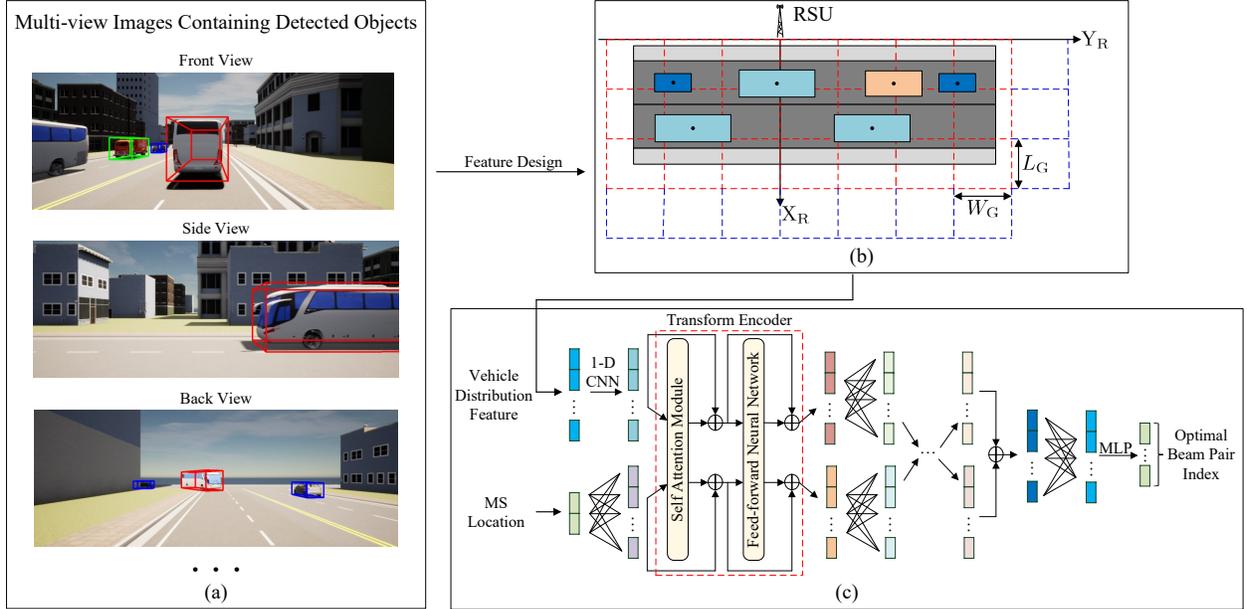


Fig. 2. The diagram of the proposed VBALA.

$X_R$ -axis,  $Y_R$ -axis, and  $Z_R$ -axis, where the origin,  $Y_R$ -axis and  $X_R$ -axis are the RSU location, the right direction of traffic lane and the vertical direction of traffic lane respectively, while  $Z_R$ -axis is upward and is perpendicular to the lane plane.

Note that, for an object, e.g., a vehicle or camera, we define the angle between the object orientation and the  $Y$ -axis as the azimuth of the object in that specific coordinate system, as shown in Fig. 1. The azimuth is positive for clockwise rotation. For clarity, we utilize the subscript to distinguish the reference coordinate system of a coordinate value or an azimuth. For instance,  $(x_M, y_M, z_M)$ ,  $(x_{C,i}, y_{C,i}, z_{C,i})$  and  $(x_R, y_R, z_R)$  denote the coordinate value under the MCS, the  $i$ th CCS, and the RCS, respectively. Moreover,  $\theta_M$ ,  $\theta_{C,i}$  and  $\theta_R$  denote the azimuth under MCS, the  $i$ th CCS, and RCS, respectively.

The locations and the azimuths of the  $i$ th camera are fixed in MCS and are then denoted by  $(x_M^i, y_M^i, z_M^i)$  and  $\theta_M^i$ ,  $i = 1, 2, \dots, C$  respectively. Without loss of generality, we assume the orientations  $\theta_M^i$  and the horizontal field of view (HFOV) of the equipped cameras can cover the horizontal view of 360 degrees around the MS, and the shooting scope of each camera does not overlap with each other. Assume there are  $O_i$  vehicles in the  $i$ th image,  $i = 1, 2, \dots, C$ , as shown in Fig. 2(a).<sup>1</sup> MS can apply 3D detection [27] to detect the  $j$ th vehicles in the  $i$ th image,

<sup>1</sup>The value of  $O_i$  could be different for different camera view

denoted as the  $(i, j)$ th vehicle,  $j = 1, 2, \dots, O_i$ . Moreover, 3D detection technique [27] can also be used to obtain the length  $l_{i,j}$ , width  $w_{i,j}$ , height  $h_{i,j}$ , the center location  $(x_{C,i}^j, y_{C,i}^j, z_{C,i}^j)$  and the azimuth  $\theta_{C,i}^j$  of the  $(i, j)$ th vehicle. All these detected vehicles' size/location/orientation parameters will be utilized latter to design the input feature of the DNN for the optimal beam alignment.

According to whether the MS knows its location, we separately discuss the following two cases:

#### A. MS Location Is Known by Itself

We propose the VBALA that contains following three key steps:

- 1) Transform the coordinates and azimuths of the vehicles from CCSs to the RCS.
- 2) Perform grid quantization for the coordinates and azimuths of the vehicles under RCS to obtain a vehicle distribution feature (VDF).
- 3) Design a VDF based beam alignment DNN (VDBAN) to infer the optimal beam pair from the VDF and MS location.

Specifically, we obtain the coordinates  $(x_M^{i,j}, y_M^{i,j}, z_M^{i,j})$  and azimuth  $\theta_M^{i,j}$  of the  $(i, j)$ th vehicle through the coordinate transformation between the MCS and the  $i$ th CCS:

$$\begin{bmatrix} x_M^{i,j} \\ y_M^{i,j} \\ z_M^{i,j} \end{bmatrix} = \begin{bmatrix} \cos(\theta_M^i) & \sin(\theta_M^i) & 0 \\ -\sin(\theta_M^i) & \cos(\theta_M^i) & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_{C,i}^j \\ y_{C,i}^j \\ z_{C,i}^j \end{bmatrix} + \begin{bmatrix} x_M^i \\ y_M^i \\ z_M^i \end{bmatrix}, \quad (5)$$

$$\theta_M^{i,j} = \theta_M^i + \theta_{C,i}^j, j = 1, 2, \dots, O_i, i = 1, 2, \dots, C.$$

Since the axes of MCS and RCS are parallel to each other, we can further get the coordinates  $(x_R^{i,j}, y_R^{i,j}, z_R^{i,j})$  and azimuth  $\theta_R^{i,j}$  of the  $(i, j)$ th vehicle from  $(x_R^{i,j}, y_R^{i,j}, z_R^{i,j}) = (x_M^{i,j} + x_R^{MS}, y_M^{i,j} + y_R^{MS}, z_M^{i,j} + z_R^{MS})$  and  $\theta_R^{i,j} = \theta_M^{i,j}$ , where  $(x_R^{MS}, y_R^{MS}, z_R^{MS})$  are the MS's coordinate in RCS.

Let us then divide  $X_R$ - $Y_R$  plane of RCS into many grids with equal size, while denote the length and width of each grid as  $L_G$  and  $W_G$ , respectively. We obtain those grids that intersect with the traffic lane, as shown in Fig. 2(b), and denote the number of the obtained grids as  $G$ . The area of the  $g$ th grid can be expressed as  $\mathcal{G}_g = \{(x_R, y_R, z_R) \mid i_g^X L_G \leq x_R < (i_g^X + 1)L_G, i_g^Y W_G \leq y_R < (i_g^Y + 1)W_G, z_R = 0\}$ ,  $g = 1, 2, \dots, G$ , where  $(i_g^X L_G, i_g^Y W_G, 0)$  is a vertex of the grid and both  $i_g^X$  and  $i_g^Y$  are integers. Denote  $\mathcal{V}_g$  as the index set of the vehicles whose  $X_R$ - $Y_R$  plane locations are contained in  $\mathcal{G}_g$ , with  $\mathcal{V}_g = \{(i, j) \mid (x_R^{i,j}, y_R^{i,j}, 0) \in \mathcal{G}_g, j = 1, 2, \dots, O_i, i =$

$1, 2, \dots, C\}$ . Denote the maximum length, width and height of the vehicles in  $\mathcal{V}_g$  as  $l_{\max,g}$ ,  $w_{\max,g}$  and  $h_{\max,g}$ , respectively. We then normalize  $l_{\max,g}$ ,  $w_{\max,g}$  and  $h_{\max,g}$  by the maximum vehicle length  $L_{\max}$ , width  $W_{\max}$  and height  $H_{\max}$  of all possible vehicle types respectively, and obtain  $l_{\max,g}^N = \frac{l_{\max,g}}{L_{\max}}$ ,  $w_{\max,g}^N = \frac{w_{\max,g}}{W_{\max}}$  and  $h_{\max,g}^N = \frac{h_{\max,g}}{H_{\max}}$ . The average value  $\theta_R^g$  of the azimuths of the vehicles in  $\mathcal{V}_g$  can be computed as  $\theta_R^g = \frac{1}{\text{Card}(\mathcal{V}_g)} \sum_{(i,j) \in \mathcal{V}_g} \theta_R^{i,j}$ .

Then, we design a VDF that can represent the vehicle distribution relative to RSU. The VDF is defined as a  $G \times 4$  dimensional matrix  $\mathbf{F} \in \mathbb{R}^{G \times 4}$ , and the  $g$ th row of  $\mathbf{F}$  is set as  $[l_{\max,g}^N, w_{\max,g}^N, h_{\max,g}^N, \theta_R^g]$ . When one grid does not contain vehicles, the corresponding row of  $\mathbf{F}$  will be set as a zero vector. Thus, the VDF  $\mathbf{F}$  can reflect the distribution of MS's surrounding vehicles.<sup>2</sup>

Next, we use transformer model [28] and design a VDBAN that can fuse the VDF  $\mathbf{F}$  and the known MS's location  $(x_R^{\text{MS}}, y_R^{\text{MS}}, z_R^{\text{MS}})$ , as shown in the Fig. 2(c). The self-attention module of transformer can conveniently perform the information exchange between data with different modalities [29]-[30]. We here adopt the TransFuser model architecture [31]. Since  $z_R^{\text{MS}}$  is constant, the input features of VDBAN include the plane coordinates  $[x_R^{\text{MS}}, y_R^{\text{MS}}]$  and the VDF  $\mathbf{F}$ . The matrix  $\mathbf{F}$  will be input into several 1-dimensional (1-D) convolution layers and fully connected (FC) layers to obtain the vector  $\mathbf{f} \in \mathbb{R}^{1 \times D}$ , and  $[x_M, y_M]^T$  will also be fed into FC layers to obtain the vector  $\mathbf{u} \in \mathbb{R}^{1 \times D}$ .

Then, the self-attention module will utilize the trainable weight matrices  $\mathbf{W}_Q \in \mathbb{R}^{D \times d}$ ,  $\mathbf{W}_K \in \mathbb{R}^{D \times d}$  and  $\mathbf{W}_V \in \mathbb{R}^{D \times d}$  to generate the *queries*, the *keys* and the *values* for  $\mathbf{u}$  and  $\mathbf{f}$  respectively, whereas the *queries*, the *keys* and the *values* of  $\mathbf{u}$  are then given by  $\mathbf{q}_u = \mathbf{u}\mathbf{W}_Q$ ,  $\mathbf{k}_u = \mathbf{u}\mathbf{W}_K$  and  $\mathbf{v}_u = \mathbf{u}\mathbf{W}_V$  respectively. Similarly,  $\mathbf{q}_f = \mathbf{f}\mathbf{W}_Q$ ,  $\mathbf{k}_f = \mathbf{f}\mathbf{W}_K$  and  $\mathbf{v}_f = \mathbf{f}\mathbf{W}_V$  are the *queries*, the *keys* and the *values* of  $\mathbf{f}$  respectively. The dot products between each *query* and all *keys* are input into the Softmax function to determine the weights for all the *values*, while the weighted sum of the *values* is used as the layer output corresponding to the *query*. Specifically, the output

<sup>2</sup>Although the grid quantization operation only uses one cubic block to represent all the vehicles contained in a grid, one will see in the later simulation that the proposed VDF can still help provide better beam alignment than the existing methods.

of self-attention module of  $\mathbf{u}$  and  $\mathbf{f}$  are denoted by  $\mathbf{u}'$  and  $\mathbf{f}'$  respectively, where

$$\begin{aligned} \begin{bmatrix} \mathbf{u}' \\ \mathbf{f}' \end{bmatrix} &= \text{Softmax}\left(\frac{1}{\sqrt{d}}\mathbf{Q}\mathbf{K}^T\right)\mathbf{V}, \\ \mathbf{Q} &= \begin{bmatrix} \mathbf{q}_u \\ \mathbf{q}_f \end{bmatrix}, \quad \mathbf{K} = \begin{bmatrix} \mathbf{k}_u \\ \mathbf{k}_f \end{bmatrix}, \quad \mathbf{V} = \begin{bmatrix} \mathbf{v}_u \\ \mathbf{v}_f \end{bmatrix}. \end{aligned} \quad (6)$$

We next utilize the *multi-head* attention mechanism [28] to enhance the representation ability of Transformer. For *multi-head* attention, denote  $h$  as the number of heads. We can utilize  $h$  groups of the attention weight matrices ( $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V$ ) to obtain different self-attention module outputs  $(\mathbf{u}'_1, \mathbf{f}'_1), (\mathbf{u}'_2, \mathbf{f}'_2), \dots, (\mathbf{u}'_h, \mathbf{f}'_h)$  according to equation (6). All the outputs of multi-head are concatenated and are mapped to the  $D$ -dimensional vector  $\mathbf{u}_o$  and  $\mathbf{v}_o$  by a linear mapping  $\mathbf{W}_O \in \mathbb{R}^{d*h \times D}$ , i.e.,

$$\begin{bmatrix} \mathbf{u}_o \\ \mathbf{f}_o \end{bmatrix} = \begin{bmatrix} \text{Concat}(\mathbf{u}'_1, \mathbf{u}'_2, \dots, \mathbf{u}'_h) \\ \text{Concat}(\mathbf{f}'_1, \mathbf{f}'_2, \dots, \mathbf{f}'_h) \end{bmatrix} \mathbf{W}_O. \quad (7)$$

The self-attention module and the feed-forward network formulate a transformer encoder and both integrated with the residual connection [32]. The outputs of the transformer encoder are two  $D$ -dimensional vectors that are obtained by feeding  $\mathbf{u} + \mathbf{u}_o$  and  $\mathbf{f} + \mathbf{f}_o$  into the same feed-forward network with two FC layers respectively. The obtained outputs of the transformer encoder will then be repeatedly fed into the FC layers and the transformer encoder module. The two outputs from the final transformer encoder module will be added up and then be input into a multi-layer perception (MLP) network to generate the output of VDBAN. Denote the set of all possible beam pairs from codebooks  $\mathcal{W}_B$  and  $\mathcal{W}_U$  as  $\mathcal{W}_P = \{(\mathbf{w}_{B,1}, \mathbf{w}_{U,1}), (\mathbf{w}_{B,1}, \mathbf{w}_{U,2}), \dots, (\mathbf{w}_{B,N_B^{\text{CB}}}, \mathbf{w}_{U,N_U^{\text{CB}}})\}$ , where  $\text{Card}(\mathcal{W}_P) = N_B^{\text{CB}} N_U^{\text{CB}}$ . The output of VDBAN is the index of the optimal beam pair in set  $\mathcal{W}_P$ . With the predicted optimal beam pair index from VDBAN, the optimal transmit and receive beam can be indexed in  $\mathcal{W}_B$  and  $\mathcal{W}_U$  respectively. Note that, the proposed beam alignment method is performed at MS and the index of the optimal transmit beam is fed back to RSU.

### B. MS Location Is Not Known

When the MS's location information cannot be obtained accurately and immediately, the VDBAN cannot be applied or would present severe performance loss. In this case, we consider to exploit the background information in the images, since the appearance and the distribution

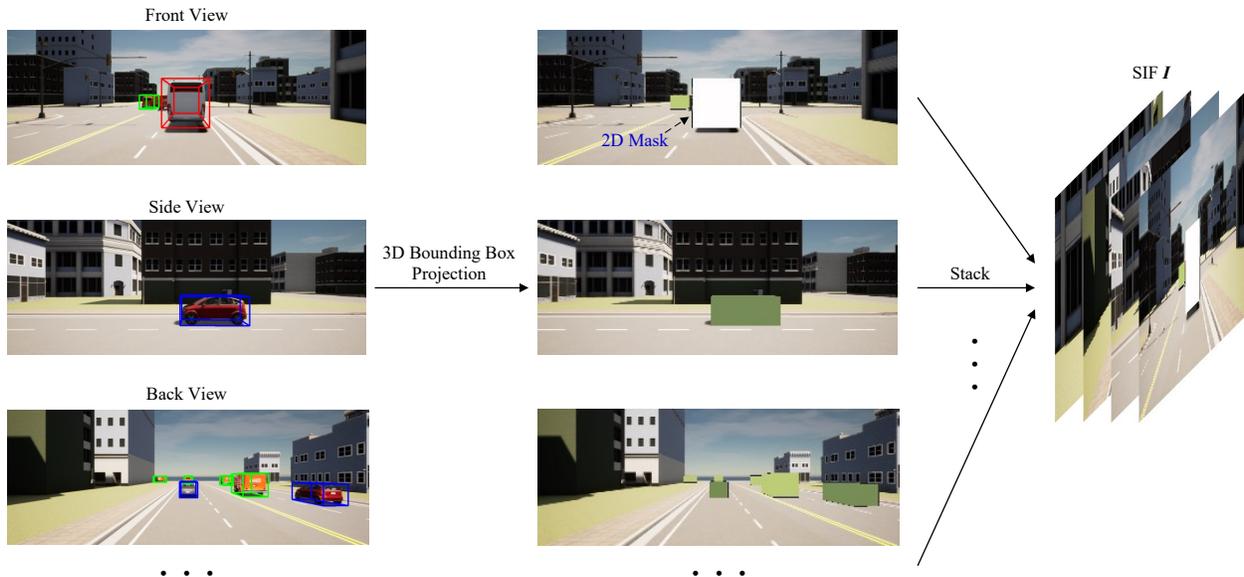


Fig. 3. The design of SIF for VBALU.

information of surrounding buildings can effectively reflect the MS location information [18]. Furthermore, the distribution of the pixels of vehicles in the images, i.e., the foreground information of images, can also intuitively reflect the vehicle distribution relative to the MS. For example, the orange van shown in the left part of the image taken at front view in Fig. 2(a) can approximately indicate the existence of a van on the left-front of the MS. Thus, the images can indicate both the MS's location information and the scattering characteristics of surrounding vehicles, and thereby can be used to infer the best beam pair without the MS's location.

However, the size of pixel area of the vehicle cannot represent the vehicle's actual size and distance from MS. For instance, the vehicle with small size can still look very large in the image when the vehicle is close to the camera. This is the inherent defect of 2D visual perception of the ordinary camera, but can be well handled by the 3D perception of Lidar or RGB-D camera. Moreover, compared with the size/location/orientation of vehicles, the color information of pixels of the vehicles is unimportant and is redundant for beam alignment, as the vehicle color may have almost no impact on the channel propagation but will increase the learning difficulty of DNN.

Hence, we design a scene image feature (SIF) by converting the vehicle color information to the vehicle size information. Specifically, we replace the three color components of the pixels of each vehicle with the length, width and height of the vehicle. For the  $(i, j)$ th vehicle, we

utilize  $(-255\frac{l_{i,j}}{L_{\max}}, -255\frac{w_{i,j}}{W_{\max}}, -255\frac{h_{i,j}}{H_{\max}})$  as the RGB channel values for all pixels of the vehicle, where the minus sign is used to distinguish the pixels of vehicles from the pixels of the scene background. We estimate the 3D bounding box of the vehicle by 3D detection. As Fig. 3 shows, we project the 3D bounding box onto the image to generate a compact 2D mask. The pixels of the vehicle are considered to be the pixels contained in the 2D mask. In fact, the vehicles pixels can be extracted more precisely by the foreground or background segmentation method [33]. However, the 2D mask projected by 3D bounding box is utilized here for simplicity. Although the adopted 2D mask is slightly coarse to cover the vehicle pixels, the effectiveness of the VBALU can still be demonstrated through the simulation.

When the vehicles in all the  $C$  images are masked with the size information, we concatenate the  $C$  images into a  $3C$ -channel image to formulate the SIF  $I$ . Specifically, the SIF  $I \in \mathbb{R}^{f_H \times f_W \times 3C}$  is defined as a 3D matrix, where  $f_H$  and  $f_W$  are the height and the width of the image respectively. The  $i$ th image with 2D masks is used as  $I_{[:,:(3c-2):3c]}$ ,  $i = 1, 2, \dots, C$ . Then, we design an SIF based beam alignment DNN (SIBAN) that adopts the widely-used ResNet architecture, such as the well-known ResNet-34 or ResNet-50 architecture [32], to predict the optimal beam pair from SIF. The input of the SIBAN is  $I$  and the output is the same as VDBAN.

**Remark:** It is worth noting that though VBALA requires the MS's location that is unnecessary for VBALU, VBALA has better environmental adaptability and generalization than VBALU. Specifically, the beam alignment performance of VBALA is almost unaffected by the surrounding environmental buildings, as the main scatters are the vehicles in traffic lane. However, since the VBALU need to infer MS's location information from the background information of images, the beam alignment performance of VBALU will be related to the spatial/color characteristics of environmental buildings.

#### IV. VISION BASED BCT PREDICTION METHOD

The accurate acquisition of BCT, i.e., the duration of the optimal beam pair, is crucial to improve the transmission rate. As Fig. 4 shows, the estimated BCT  $T_L$  that is longer than the actual BCT  $T_A$  can cause severe beam misalignment during some time periods, such as the time period  $[T_A, T_L]$  and thereby results in lower transmission rate than that can be achieved by the actual BCT. When the estimated BCT  $T_S$  is shorter than the actual BCT  $T_A$ , the beam alignment will be frequently performed, as Fig. 4 shows the beam alignment is repeated for 3 times during

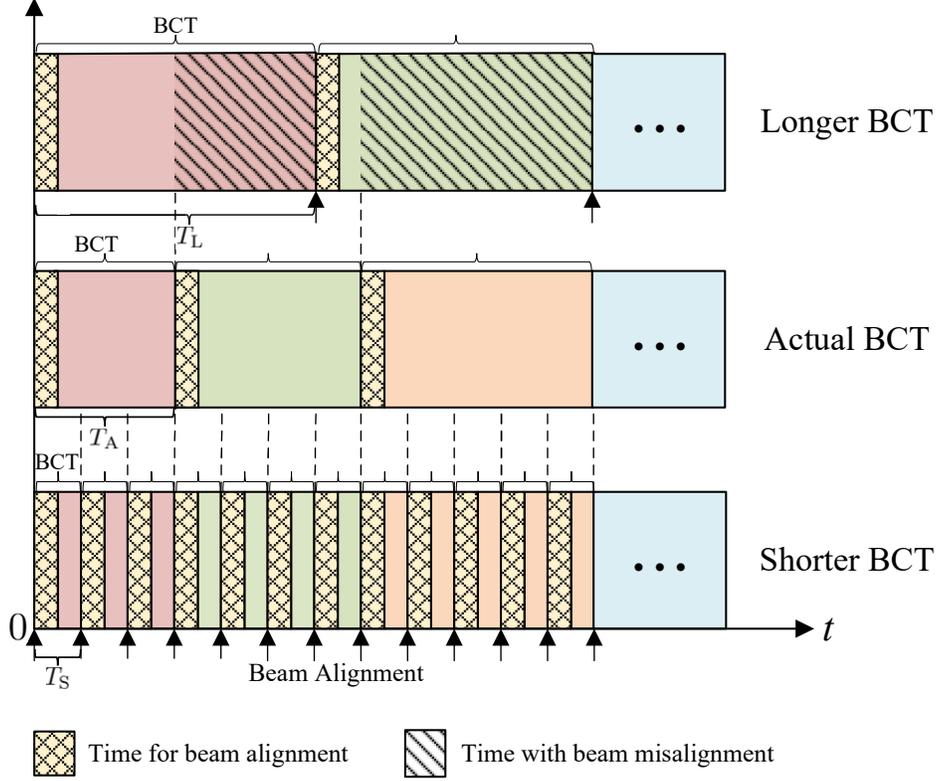


Fig. 4. Trade-off between the length of BCT and the transmission rate.

the actual BCT period  $[0, T_A]$ . Thus, the transmission rate will also decrease due to the large time overhead for beam alignment.

Specifically, the achievable transmission rate is affected by two factors: the beam alignment accuracy and the beam alignment time. Denote the total time for communications as  $T_{\text{total}}$ , and denote the total time for beam alignment as  $T_{\text{align}}$ . Moreover, the average transmission rate achieved by the beam alignment for the time  $T_{\text{total}} - T_{\text{align}}$  without beam alignment overhead is denoted as  $R_{\text{ave}}$ , and  $R_{\text{ave}}$  is positively correlated with the beam alignment accuracy. Thus, the actual average transmission rate considering the time overhead of beam alignment can be expressed as  $(1 - \frac{T_{\text{align}}}{T_{\text{total}}})R_{\text{ave}}$ . When the estimated BCT is longer than the actual BCT,  $T_{\text{align}}$  can become shorter, but  $R_{\text{ave}}$  may severely decrease due to the beam misalignment. When the estimated BCT is shorter than the actual BCT, the  $R_{\text{ave}}$  may be optimal, but the  $T_{\text{align}}$  will increase to degrade the actual average transmission rate. Thus, the accurate BCT is expected to obtain to realize the optimal trade-off between  $R_{\text{ave}}$  and  $T_{\text{align}}$ .

The BCT is mainly affected by the location, size, moving direction and speeds of the MS

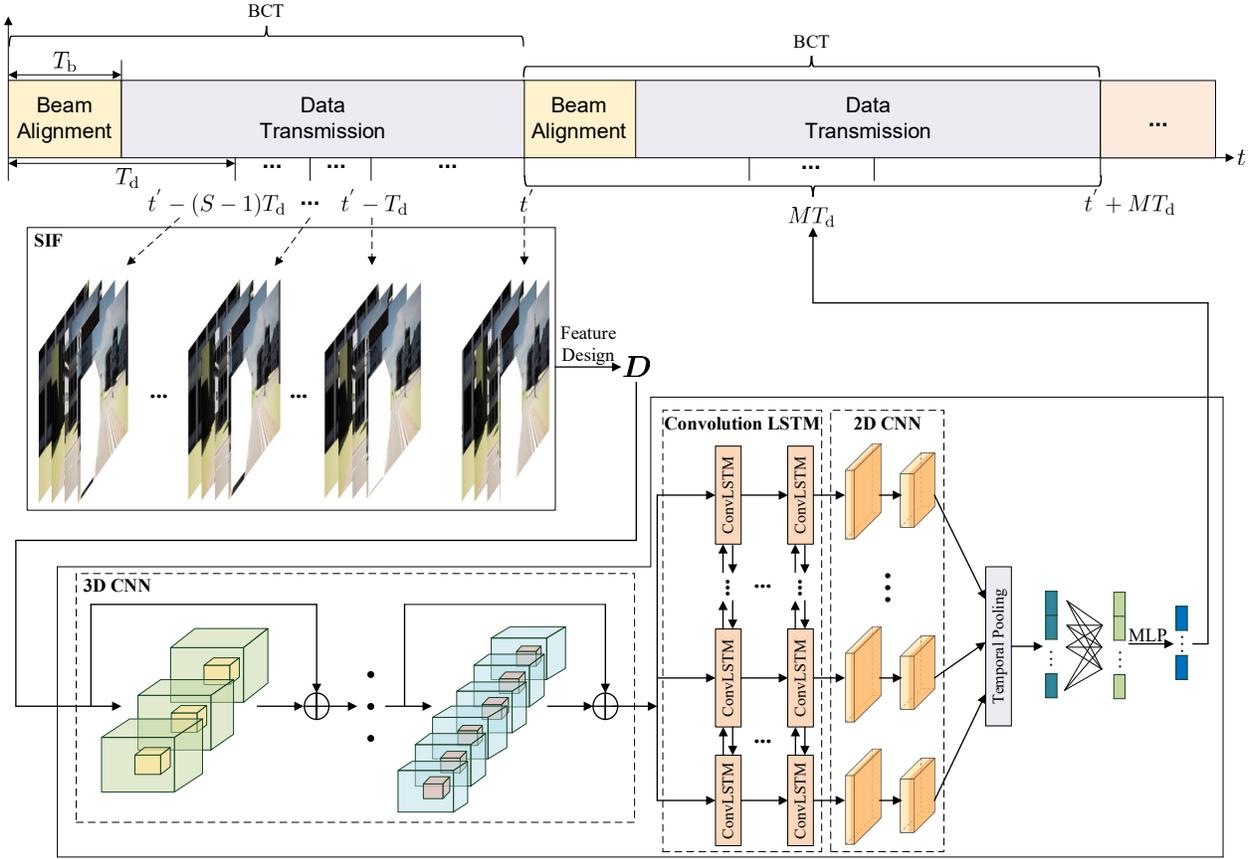


Fig. 5. The diagram of the proposed VPBCT.

and the surrounding scatters. Conventionally, the BCT can only be estimated under some LOS scenarios ordinarily, such as the high speed train communication and unmanned aerial vehicle communication, from the MS velocity and the size of beam coverage area [34]-[35]. For the complex environment with dynamic scatters and blockages, the traditional beam alignment strategy generally adopts the fixed BCT that is not guaranteed to be accurate [13]-[15],[36]. As the consecutively perceived scene images at MS can effectively present the exact spatial information and moving characteristics of the environmental objects, we use DNN to predict the accurate BCT from the scene images.

The diagram of the proposed VPBCT method is shown in Fig. 5. Denote  $T_d$  as the shooting interval of the camera, i.e., the reciprocal of the camera's frame rate, and denote  $T_b$  as the time for beam alignment during one BCT. During each BCT, the communication system will first perform beam alignment to estimate the optimal beam pair and then utilize the estimated beam pair for data transmission. Since the actual BCT is a continuous variable, we discretize BCT

with duration  $T_d$  to reduce the learning difficulty of the DNN. Without loss of generality, we express the BCT of the optimal beam pair at moment  $t'$  as  $MT_d$ , where  $M$  is an integer with  $M \geq 1$ . Thus, the minimum BCT that may be adopted by system is  $T_d$ , and the time  $T_b$  for beam alignment is assumed to be less than the minimum BCT, i.e.,  $T_b < T_d$ , for simplicity.

We design a scene image based BCT prediction DNN (SIBPN) that adopts the 3D CNN integrated with the ConvLSTM (3D CNN-LSTM) model [37] to predict  $M$ . As shown in Fig. 5, the SIBPN is formed by connecting 3D CNN, ConvLSTM and 2D CNN module in series. The 3D CNN module includes several 3D CNN layers with residual connection, the ConvLSTM module consists of the bidirectional ConvLSTM layers, and the 2D CNN module consists of the ordinary 2D CNNs. The ConvLSTM is designed by replacing the fully connected operators of the traditional LSTM with the convolution operator to obtain the spatial information in the images [38].

We utilize the images that are taken by the  $C$  cameras at the MS to generate the input feature  $D$  of the SIBPN. Specifically,  $D \in \mathbb{R}^{S \times f_H \times f_W \times 3C}$  is a 4-dimensional matrix, and its  $p$ th row is the SIF generated by the  $C$  images taken at moment  $t' - (S - p)T_d$ ,  $p = 1, 2, \dots, S$ . The input feature  $D$  is fed into the 3D CNN module to obtain a 4D tensor  $D' \in \mathbb{R}^{S' \times f'_H \times f'_W \times 3C'}$ . Then,  $D'$  is split into a sequential image data with  $S'$  time steps, i.e.,  $D'_{[1, \dots, :]}, D'_{[2, \dots, :]}, \dots, D'_{[S', \dots, :]}$ . The sequential image data is fed into the ConvLSTM module to extract both the temporal and spatial characteristics of the environmental objects. The output of ConvLSTM of each time step will be fed into a 2D CNN respectively for dimension reduction. Finally, the average of the output vectors of all 2D CNNs are calculated for temporal pooling, and this average vector is fed into several FC layers to produce the output of the SIBPN. Since  $M$  is assumed as an integer, the BCT prediction problem can be regarded as a classification problem. We here assume a maximum possible length of BCT as  $M_{\max}T_d$ . Thus,  $M$  must be in the set  $\mathcal{B} = \{1, 2, \dots, M_{\max}\}$ , and the output of the SIBPN is the index of the accurate BCT in set  $\mathcal{B}$ .

Once  $M$  is predicted by SIBPN, the communication system will implement beam alignment during the time period  $[t', t' + T_b]$ , and then keeps on using the optimal beam pair obtained from beam alignment during the time period  $[t' + T_b, t' + MT_d]$ . The next BCT prediction as well as the beam alignment will be performed at moment  $t' + MT_d$ .

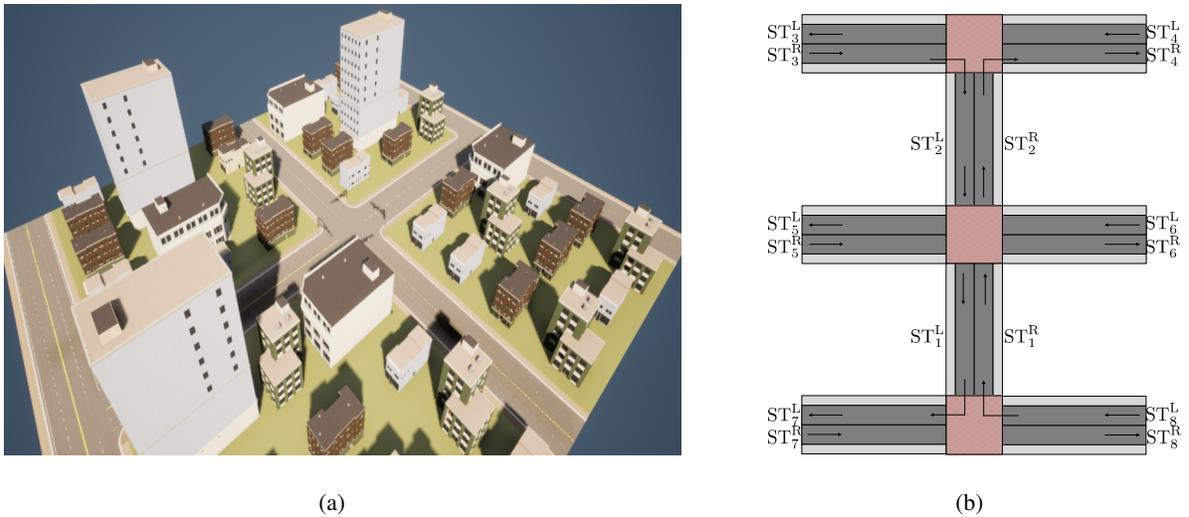


Fig. 6. (a) The constructed communication environment in CARLA. (b) The eight streets of the communication environment.

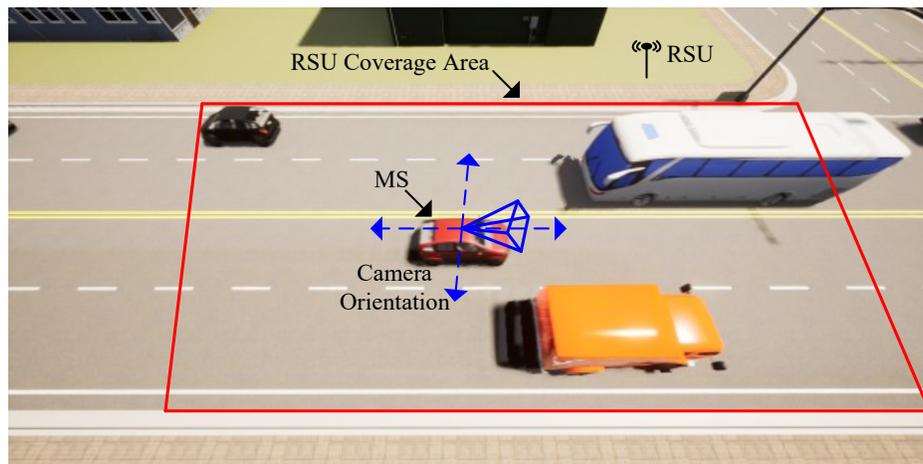


Fig. 7. The simulated RSU coverage area.

## V. SIMULATION RESULTS

In this section, we verify the performance of VBALAA, VBALAU and VPBCT in a simulated traffic V2X environment. The images and channels are collected to construct the datasets of VDBAN, SIBAN and SIBPN.

### A. Simulation Setup

1) *Environment Generation*: We utilize the CARLA [39], an autonomous driving simulation platform, to simulate the environment as well as vehicles. The CARLA could support to equip

TABLE III  
VEHICLE SIZES FOR SIMULATION

Type	Length/m	Width/m	Height/m
Car	3.71	1.79	1.55
Van	5.20	2.61	2.47
Bus	11.08	3.25	3.33

the vehicle with many kinds of sensors, such as camera, LIDAR and inertial measurement unit (IMU), etc., to obtain abundant sensory measurement data. Fig. 6(a) shows the constructed 3D communication environment model in CARLA. The environment model has eight streets. Denote the left and the right lane of the  $r$ th street as  $ST_r^L$  and  $ST_r^R$  respectively,  $r = 1, 2, \dots, 8$ , as shown in Fig. 6(b). The specified driving direction for each lane is also shown in Fig. 6(b).

The RSU is placed at the side of  $ST_1^L$ . The concerned RSU's coverage area is part of  $ST_1^L$  and  $ST_1^R$ , which is 30 meters long and 15 meters wide, as shown in the Fig. 7. Three vehicle types including *car type*, *van type* and *bus type* are adopted, whose sizes are listed in TABLE III. To simulate traffic scenario, we first randomly generate 15 vehicles in  $ST_1^L$ ,  $ST_2^L$ ,  $ST_3^R$  and also randomly generate 15 vehicles in  $ST_1^R$ ,  $ST_2^R$ ,  $ST_8^L$  for the vehicle initialization. Here, we will give priority to generate a vehicle with the *car type* in  $ST_1^R$  and set the vehicle as MS. The types of all generated vehicles and the colors of the vehicles belonging to *car type* are randomly determined. Then, we utilize the Simulation of Urban MObility (SUMO) [40], a traffic simulation software, to control the speed and moving trajectory of all the vehicles by the co-simulation interface of the CARLA.

2) *Scene Image Generation*: The locations of all  $C$  cameras are set at 0.5m above the roof center of MS, and the  $C$  cameras have different orientations. When the MS is running through the RSU's coverage area, it keeps on taking images with the interval  $T_d$  until it leaves the area. The  $C$  images taken at each moment form an image set, and the MS can collect a sequence of image sets within the RSU's coverage area.

3) *Channel Generation*: We adopt the Wireless Insite [41], a ray tracing software, to simulate the channel. Note that the channel produced by ray tracing technique can have consistent spatial correlation with the communication environment. The setup parameters of Wireless Insite are listed as TABLE IV. The ULA equipped at BS and MS are both set to be parallel with  $ST_1^L$ . The ULA of MS is equipped at 0.05m above the roof center of MS. The ULA of RSU is

TABLE IV  
CRITICAL PARAMETERS OF WIRELESS INSITE FOR RAY TRACING

Parameter	Value
Carrier Frequency	28 GHz
Propagation Model	X3D
Building Material	Concrete
Vehicle Material	Metal
Maximum Number of Reflections	6
Maximum Number of Diffractions	1
Maximum Paths Per Receiver Point	5

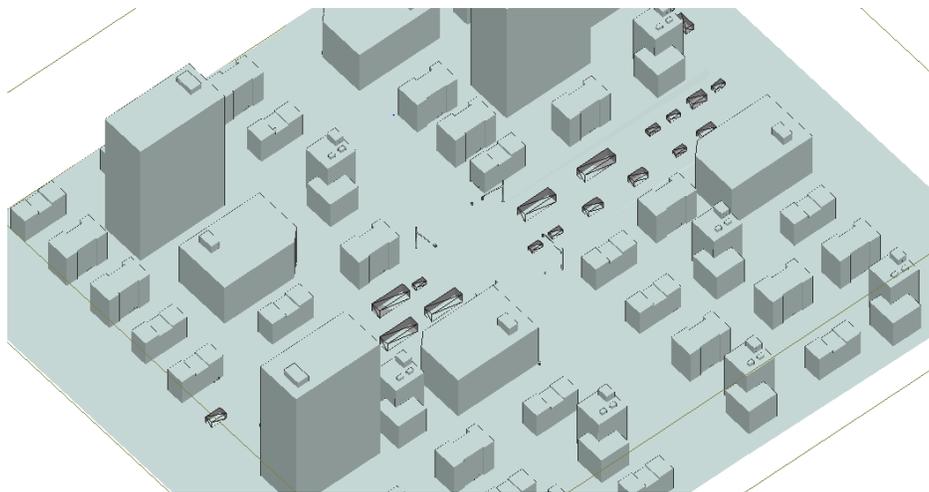


Fig. 8. The synchronization simulation in Wireless Insite.

equipped at 3m above the ground. For each moment that MS take images, we will synchronize the environment and all vehicles of the CARLA into Wireless Insite. Specifically, as Fig. 8 shows, we ensure the environment model and the sizes/locations/orientations of vehicles in Wireless Insite are exactly the same as that in CARLA to produce the corresponding channel. Thus, the produced channel can match the image set taken at the same moment.

To reduce the computation overhead of channel generation, we convert the environment model and the vehicle models of CARLA into the simple cube-like models by ignoring the surface details of the models of CARLA, and utilize the cube-like models in Wireless Insite. Since the small surface details of CARLA models can only have a slight impact on the channel, this simplification for CARLA models will not affect the reliability of simulation.

4) *DNN Dataset Generation*: According to Section V.A-1) and Section V.A-2), the types and initial locations of generated vehicles are different with different vehicle initializations, and thus the image set sequences collected by MS will be different from each other. To obtain the datasets of VDBAN/SIBAN/SIBPN, we implement  $Q$  different vehicle initializations and thereby acquire  $Q$  different image set sequences. Denote the number of image sets contained in the  $q$ th image set sequence as  $S_q$ , and denote the  $r$ th image set of the  $q$ th image set sequence as  $\mathcal{I}_{q,r}$ ,  $r = 1, 2, \dots, S_q$ ,  $q = 1, 2, \dots, Q$ . The channel corresponding to  $\mathcal{I}_{q,r}$  will be generated according to Section V.A-3), and the optimal beam pair index  $b_{q,r}^{\text{opt}}$  can be calculated from the channel. We also record the MS's location each shooting moment and denote the MS location as  $\mathbf{p}_{q,r}$ . Moreover, the VDF  $\mathbf{F}_{q,r}$  can be generated from  $\mathcal{I}_{q,r}$  by the 3D detection method.

Thus, we can pair all the VDFs and MS locations with the corresponding optimal beam pair indexes and construct the sample set  $\mathcal{S}_{\text{VDBAN}} = \{(\mathbf{F}_{q,r}, \mathbf{p}_{q,r}, b_{q,r}^{\text{opt}}) \mid r = 1, 2, \dots, S_q, q = 1, 2, \dots, Q\}$ . Each element in  $\mathcal{S}_{\text{VDBAN}}$  can be a training/validation/test sample of VDBAN. We divide the  $\mathcal{S}_{\text{VDBAN}}$  into the training set  $\mathcal{S}_{\text{VDBAN}}^{\mathcal{Q}_T}$ , the validation set  $\mathcal{S}_{\text{VDBAN}}^{\mathcal{Q}_V}$  and test set  $\mathcal{S}_{\text{VDBAN}}^{\mathcal{Q}_E}$  of VDBAN, where  $\mathcal{S}_{\text{VDBAN}}^{\mathcal{Q}} = \{(\mathbf{F}_{q,r}, \mathbf{p}_{q,r}, b_{q,r}^{\text{opt}}) \mid r = 1, 2, \dots, S_q, q \in \mathcal{Q}\}$  and  $\mathcal{Q} = \mathcal{Q}_T, \mathcal{Q}_V, \mathcal{Q}_E$ . Moreover, the  $\mathcal{Q}_T, \mathcal{Q}_V$  and  $\mathcal{Q}_E$  are the sets of the indices of image set sequences used for constructing training, validation and test set respectively, with  $\text{Card}(\mathcal{Q}_T) + \text{Card}(\mathcal{Q}_V) + \text{Card}(\mathcal{Q}_E) = Q$ . Denote the SIF generated from  $\mathcal{I}_{q,r}$  as  $\mathbf{I}_{q,r}$ . Similarly for the SIBAN, we can obtain the training set  $\mathcal{S}_{\text{SIBAN}}^{\mathcal{Q}_T}$ , the validation set  $\mathcal{S}_{\text{SIBAN}}^{\mathcal{Q}_V}$  and the test set  $\mathcal{S}_{\text{SIBAN}}^{\mathcal{Q}_E}$ , where  $\mathcal{S}_{\text{SIBAN}}^{\mathcal{Q}} = \{(\mathbf{I}_{q,r}, b_{q,r}^{\text{opt}}) \mid r = 1, 2, \dots, S_q, q \in \mathcal{Q}\}$  and  $\mathcal{Q} = \mathcal{Q}_T, \mathcal{Q}_V, \mathcal{Q}_E$ .

To train the SIBPN, we should first determine the BCT of the optimal beam pair  $b_{q,r}^{\text{opt}}$ . For simplicity, we assume the channel will not change during a shooting interval  $T_d$  to ensure the BCT will be the integral multiple of  $T_d$ . We compare whether the element in  $\{b_{q,r+r'}^{\text{opt}} \mid r' = 1, 2, \dots, S_q - r\}$  is the same as  $b_{q,r}^{\text{opt}}$  one by one to find the maximum duration of  $b_{q,r}^{\text{opt}}$ . Specifically, if  $b_{q,r}^{\text{opt}} = b_{q,r+1}^{\text{opt}} = \dots = b_{q,r+M_{q,r}-1}^{\text{opt}} \neq b_{q,r+M_{q,r}}^{\text{opt}}$ , then we can determine the BCT of  $b_{q,r}^{\text{opt}}$  as  $M_{q,r}T_d$ . Next,  $\mathcal{I}_{q,r-S+1}, \mathcal{I}_{q,r-S+2}, \dots, \mathcal{I}_{q,r}$  can be used to generate the SIBPN's input feature  $\mathbf{D}_{q,r}$  for predicting  $M_{q,r}$ ,  $r = S, S+1, \dots, S_q$ . For SIBPN, we pair  $\mathbf{D}_{q,r}$  and  $M_{q,r}$  to construct the training set  $\mathcal{S}_{\text{SIBPN}}^{\mathcal{Q}_T}$ , the validation set  $\mathcal{S}_{\text{SIBPN}}^{\mathcal{Q}_V}$ , and the test set  $\mathcal{S}_{\text{SIBPN}}^{\mathcal{Q}_E}$ , where  $\mathcal{S}_{\text{SIBPN}}^{\mathcal{Q}} = \{(\mathbf{D}_{q,r}, M_{q,r}) \mid r = S, S+1, \dots, S_q, q \in \mathcal{Q}\}$  and  $\mathcal{Q} = \mathcal{Q}_T, \mathcal{Q}_V, \mathcal{Q}_E$ .

5) *The Adoption of 3D Detection Method*: We adopt the *single-stage monocular 3D object detection via keypoint estimation* (SMOKE) in [27] as an example to support VBALA, VBALU and VPBCT. To train the SMOKE network, we implement  $Q_D$  different vehicle initializations,

and let MS take images in  $ST_1^R$  and  $ST_2^R$  for each vehicle initialization. All the images will be labeled with the 3D bounding boxes of the contained vehicles to construct the samples for the SMOKE network.

### B. Performance Metrics

For VBALA and VBALU, we utilize the achievable transmission rate ratio (ATRR), i.e., the ratio between the average transmission rate achieved by the selected beam pair and the optimal average transmission rate as the performance metric. Denote the  $k$ th subcarrier channel matrix corresponding to the image set  $\mathcal{I}_{q,r}$  as  $\mathbf{H}_{q,r,k}$ ,  $r = 1, 2, \dots, S_q$ ,  $q = 1, 2, \dots, Q$ ,  $k = 1, 2, \dots, K$ , and denote the shooting moment that MS takes  $\mathcal{I}_{q,r}$  as  $t_{q,r}$ . Since  $\mathbf{H}_{q,r,k}$  is assumed unchanged during an interval  $T_d$ , the optimal transmission rate during the time period  $[t_{q,r}, t_{q,r} + T_d]$  is given by

$$R(\mathbf{w}_{B,q,r}^{\text{opt}}, \mathbf{w}_{U,q,r}^{\text{opt}}) = \frac{1}{K} \sum_{k=1}^K \log_2 \left( 1 + \frac{P_k}{\sigma^2} |(\mathbf{w}_{U,q,r}^{\text{opt}})^H \mathbf{H}_{q,r,k} \mathbf{w}_{B,q,r}^{\text{opt}}|^2 \right), \quad (8)$$

where  $\mathbf{w}_{B,q,r}^{\text{opt}}$  and  $\mathbf{w}_{U,q,r}^{\text{opt}}$  are the transmit and receive beam corresponding to the beam pair index  $b_{q,r}^{\text{opt}}$  respectively. Then, we assume the communication system will perform beam alignment at moment  $t_{q,r}$ ,  $r = 1, 2, \dots, S_q$ ,  $q = 1, 2, \dots, Q$ . Denote the transmit and receive beam selected by VBALA or VBALU at shooting moment  $t_{q,r}$  as  $\mathbf{w}_{B,q,r}$  and  $\mathbf{w}_{U,q,r}$  respectively. The ATRR for VBALA or VBALU can be expressed

$$\text{ATRR}_s^{\mathcal{Q}} = \frac{\sum_{q \in \mathcal{Q}} \sum_{r=1}^{S_q} R(\mathbf{w}_{B,q,r}, \mathbf{w}_{U,q,r})}{\sum_{q \in \mathcal{Q}} \sum_{r=1}^{S_q} R(\mathbf{w}_{B,q,r}^{\text{opt}}, \mathbf{w}_{U,q,r}^{\text{opt}})}, \quad (9)$$

where  $\mathcal{Q} = \mathcal{Q}_V, \mathcal{Q}_E$ , and  $\text{ATRR}_s^{\mathcal{Q}_V}$  and  $\text{ATRR}_s^{\mathcal{Q}_E}$  represent the ATRR of VBALAA or VBALAU on validation and test set, respectively.

For VPBCT, the BCT prediction accuracy (BCTPA) is used as the performance metric. Denote the BCT predicted from  $\mathbf{D}_{q,r}$  as  $M'_{q,r} T_d$ . The BCTPA is expressed as

$$\text{BCTPA}^{\mathcal{Q}} = \frac{\sum_{q \in \mathcal{Q}} \sum_{r=S}^{S_q} \mathbb{1}^A(M'_{q,r})}{\sum_{q \in \mathcal{Q}} (S_q - S + 1)}, \quad (10)$$

where  $\mathbb{1}^A(M'_{q,r})$  is the indicator function

$$\mathbb{1}^A(M'_{q,r}) = \begin{cases} 1, & M'_{q,r} = M_{q,r} \\ 0, & M'_{q,r} \neq M_{q,r} \end{cases}. \quad (11)$$

TABLE V  
THE PARAMETERS OF LAYERS OF SIBPN

Layer Order	Kernel/Pool Size	Strides	Filters	Units
3D Convolution	(1, 7, 7)	(1, 2, 2)	64	None
3D MaxPooling	(1, 3, 3)	(1, 2, 2)	None	None
3D Convolution	(2, 3, 3)	(1, 1, 1)	64	None
3D Convolution	(2, 3, 3)	(1, 2, 2)	128	None
3D Convolution	(2, 3, 3)	(1, 1, 1)	128	None
Bidirectional ConvLSTM	(3, 3)	(1, 1)	128	None
2D Convolution	(3, 3)	(2, 2)	64	None
2D Convolution	(3, 3)	(2, 2)	32	None
Average	None	None	None	None
FC	None	None	None	1024
FC	None	None	None	1024
FC	None	None	None	3

Moreover, there is  $\mathcal{Q} = \mathcal{Q}_V, \mathcal{Q}_E$ , and  $\text{BCTPA}^{\mathcal{Q}_V}$  and  $\text{BCTPA}^{\mathcal{Q}_E}$  represent the BCTPA of VPBCT on validation and test set, respectively. We further consider utilizing the ATRR as the performance metric to verify the improvement in transmission rate brought by BCT prediction. We assume the beam pair obtained from beam alignment is always the optimal for simplicity. Similarly, denote the transmit and receive beam adopted at shooting moment  $t_{q,r}$  according to the predicted BCTs as  $\mathbf{w}_{B,q,r}^{\text{BCT}}$  and  $\mathbf{w}_{U,q,r}^{\text{BCT}}$  respectively. Since the beam alignment will consume the time  $T_b$  in the first interval  $T_d$  of each BCT, the ATRR for VPBCT is expressed as

$$\text{ATRR}_p^{\mathcal{Q}} = \frac{\sum_{q \in \mathcal{Q}} \sum_{r=S}^{S_q} (1 - \mathbb{1}^R(t_{q,r}) \frac{T_b}{T_d}) R(\mathbf{w}_{B,q,r}^{\text{BCT}}, \mathbf{w}_{U,q,r}^{\text{BCT}})}{\sum_{q \in \mathcal{Q}} \sum_{r=S}^{S_q} R(\mathbf{w}_{B,q,r}^{\text{opt}}, \mathbf{w}_{U,q,r}^{\text{opt}})}, \quad (12)$$

where  $\mathbb{1}^R(t_{q,r})$  is the indicator function

$$\mathbb{1}^R(t_{q,r}) = \begin{cases} 1, & \text{beam alignment happens during } [t_{q,r}, t_{q,r} + T_d] \\ 0, & \text{otherwise} \end{cases}, \quad (13)$$

$\mathcal{Q} = \mathcal{Q}_V, \mathcal{Q}_E$ , and  $\text{ATRR}_p^{\mathcal{Q}_V}$  and  $\text{ATRR}_p^{\mathcal{Q}_E}$  represent the ATRR of VPBCT on validation and test set, respectively.

### C. Simulation Parameters

The numbers of antennas at BS and MS are both set as  $N_B = N_U = 64$ ;  $K$  is set as 16;  $N_B^{\text{CB}} = N_B = 64$ ;  $N_U^{\text{CB}} = N_U = 64$ ;  $\mathbf{w}_{B,b} = \mathbf{a}_r(\frac{2b-2-N_B^{\text{CB}}}{2N_B^{\text{CB}}}\pi)$ ,  $b = 1, 2, \dots, N_B^{\text{CB}}$  and  $\mathbf{w}_{U,u} =$

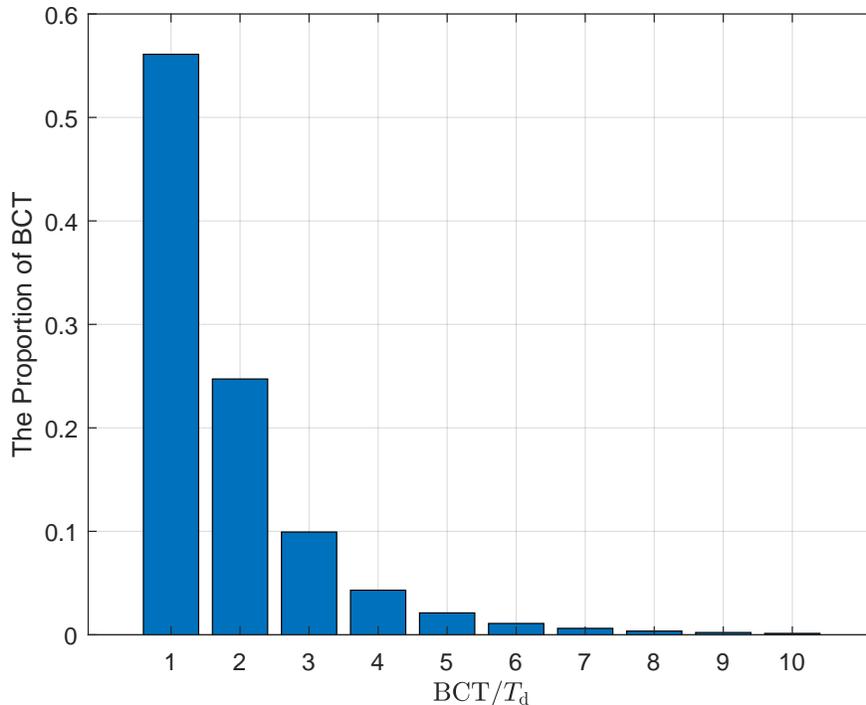


Fig. 9. The proportions of different BCTs in simulation.



Fig. 10. A case of the image taken by the camera adopted by BMBA.

$\mathbf{a}_t(\frac{2u-2-N_U^{CB}}{2N_U^{CB}}\pi)$ ,  $u = 1, 2, \dots, N_U^{CB}$ ;  $P_1 = \dots = P_K$ ;  $\frac{P_k}{K\sigma^2 \sum_{q=1}^Q S_q} \sum_{q=1}^Q \sum_{r=1}^{S_q} \sum_{k=1}^K \|\mathbf{H}_{q,r,k}\|_F^2 = 29.5\text{dB}$ ;  $C$  is set to be 4; the camera orientation  $\theta_M^1$ ,  $\theta_M^2$ ,  $\theta_M^3$  and  $\theta_M^4$  are set to be 0, 90, 180 and 270 degrees respectively, as shown in Fig. 7;  $T_d$  is set as 0.05s;  $L_G$  and  $W_G$  are set as 11.7m and 2m respectively;  $Q_D$  is set to be 50;  $Q$  is set to be 600;  $\text{Card}(\mathcal{Q}_T)$ ,  $\text{Card}(\mathcal{Q}_V)$  and  $\text{Card}(\mathcal{Q}_E)$  are 480, 60 and 60 respectively; The number of training samples, validation samples

and test samples for VDBAN and SIBAN are 27357, 3419 and 3340, respectively.

Due to the limited RSU's coverage area, not all the  $N_B N_U$  beam pairs in  $\mathcal{W}_P$  can potentially serve as the optimal beam pair. Thus, we count all different optimal beam pairs from all the used beam pairs, i.e.,  $b_{q,r}^{\text{opt}}$ ,  $r = 1, 2, \dots, S_q$ ,  $q = 1, 2, \dots, Q$ , to form a beam pair set  $\mathcal{W}'_P$ , where  $\text{Card}(\mathcal{W}'_P) = 365$ . The outputs of VDBAN and SIBAN both are set as the index of the optimal pair in set  $\mathcal{W}'_P$  instead of  $\mathcal{W}_P$ . Since  $\text{Card}(\mathcal{W}'_P)$  is much smaller than  $\text{Card}(\mathcal{W}_P)$ , the simulation overhead and practical network complexity can be reduced.

The number  $S$  of image sets used for BCT prediction is set to be 3. The adopted image resolution is  $320 \times 120$ . According to the statistics for  $M_{q,r}$ ,  $r = 1, 2, \dots, S_q$ ,  $q = 1, 2, \dots, Q$ , the proportions of different BCTs are obtained and shown in Fig. 9. It is seen that there exists serious class imbalance problem, since the proportions of short BCTs are much greater than the proportions of long BCTs. Hence, we group the BCTs to reduce the proportion difference between different BCTs. Specifically, we divide all the possible BCTS into three groups  $\{1\}$ ,  $\{2\}$ ,  $\{3, 4, \dots, M_{\max}\}$ , and set the output of SIBPN as the index of the group containing the accurate BCT. Once a group is predicted by SIBPN, the minimum value of this group will be used as the predicted BCT. Moreover, the re-sampling of the training samples is adopted to further release the class imbalance problem and enhance the generalization performance of SIBPN. Due to the grouping operation, the metric BCTPA is modified as the prediction accuracy of BCT group. When the predicted BCT group from  $D_{q,r}$  is the group that contain the true BCT  $M_{q,r}$ ,  $\mathbb{1}^A(M'_{q,r})$  is set as 1, otherwise  $\mathbb{1}^A(M'_{q,r})$  is set as 0.

For VDBAN, the VDF  $F$  is input into a 1-D convolution layer with filter number, kernel size and stride as 1, and is then input into one FC layer with node number  $D = 64$  to produce  $f$ . The MS's location is also input into one FC layer with node number 64 to produce  $u$ . Then, two transform encoder modules with  $d = 16$  and  $d = 32$  are used, and the number of heads of both the two encodes are 4. The final MLP network of VDBAN includes three FC layers with node numbers 1024, 1024 and 365 respectively. The VDBAN is trained for 60 epochs. The network structure of the SIBAN is set as the ResNet101V2 [32], and the SIBAN is trained for 20 epochs.

For the SIBPN, the parameters of the network structure are shown in the TABLE V, where the batch normalization and ReLU activation function are used for each convolutional layer. Moreover, through the residual connection, the input of the second 3D convolutional layer is added to the output of this layer, and the input of the third 3D convolutional layer is also added

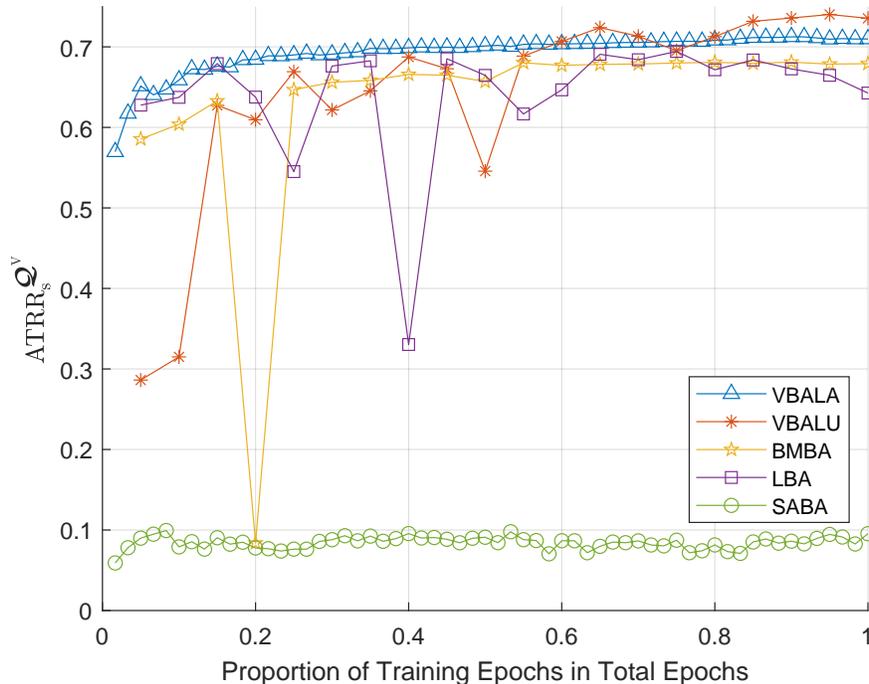


Fig. 11.  $ATTR_s^{Q^V}$  achieved by Top-1 beam pair selection with the increase of the number of training epochs.

to the output of the fourth 3D convolutional layer. The SIBPN is trained for 20 epochs.

#### D. Results and Discussions

We compare the VBALA and VBALU, with the BS's vision based multi-modal beam alignment (BMBA) in [16], the LIDAR based beam alignment (LBA) in [11] and the situational awareness based beam alignment (SABA) in [6].

The BMBA uses the images taken at BS and the MS's location as the input features to perform the beam alignment by DNN. For the camera adopted by BMBA, the plane location of the camera is set the same as that of RSU, and the height of the camera is set as 6m to monitor the entire RSU's coverage area. A case of the image taken by the camera adopted by BMBA is shown in Fig. 10. The input features and the network structure is set to be the same as that in [16] to construct the DNN used for BMBA. Moreover, for the DNN of BMBA, the feature extraction subnetwork for images is set to be the same as SIBAN, and the subnetwork for beam alignment after concatenate layer is set to have the same layer sizes and the same number of layers as VDBAN ,i.e., the seven FC layers with node numbers 64, 64, 128, 128, 1024, 1024 and 365

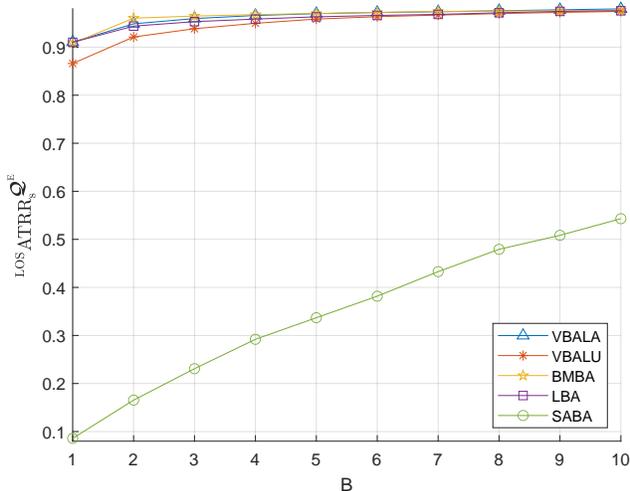


Fig. 12.  $^{LOS}ATR\bar{R}_s^E$  for Top-B beam pair selection. The number of LOS test samples are 1930, which is 58% of total test samples.

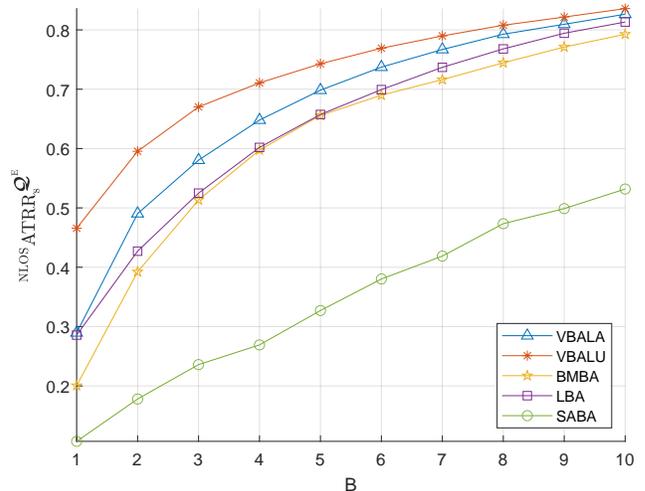


Fig. 13.  $^{NLOS}ATR\bar{R}_s^E$  for Top-B beam pair selection. The number of NLOS test samples are 1410, which is 42% of total test samples.

respectively, to keep the fairness. The number of training epochs of the DNN of BMBA is set to be consistent with SIBAN.

The LBA designs a point cloud feature (PCF) from the MS's location and the point cloud scanned at MS, and utilizes the PCF as the input feature of DNN for beam pair selection. For the DNN of LBA, the input feature is set to be the same as that in [11], while the network structure and the number of training epochs are set to be consistent with SIBAN for fairness. The SABA uses the surrounding vehicles' locations and MS location to present an environmental vehicle location feature (VLF) that is different from the VDF or SIF. Note that [6] only uses some classic machine learning methods but not the DNN for beam alignment. For the sake of fairness, we here also adopt the DNN with the same layer sizes, the same number of layers, and the same number of training epochs as VDBAN to infer the optimal beam pair from the VLF. Specifically, the DNN of SABA is constructed by seven FC layers with node numbers 64, 64, 128, 128, 1024, 1024 and 365 respectively. Moreover, the surrounding vehicle locations for VLF are obtained by the 3D detection.

As Fig. 11 shows, we firstly analyze the  $ATR\bar{R}_s^V$  achieved by Top-1 beam pair selection versus the increase of the number of training epochs for the DNNs of VBALA, VBALU, BMBA, LBA and SABA. All the methods are trained to reach the convergence, and it is seen that the VBALA and VBALU can outperform BMBA, LBA and SABA. Moreover, it is also seen that

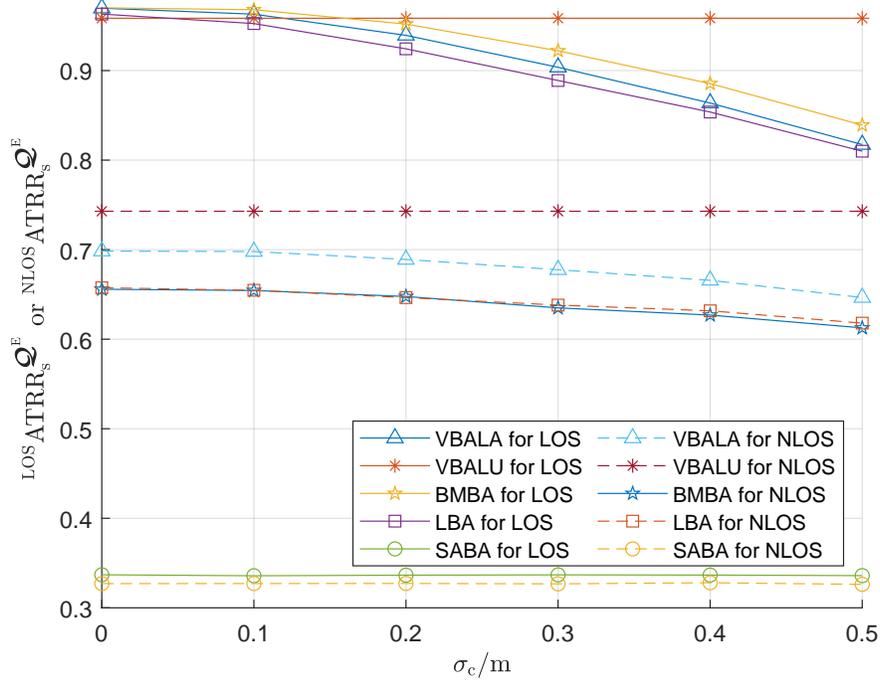


Fig. 14.  ${}^{\text{LOS}}\text{ATTR}_s^{\mathcal{Q}^E}$  and  ${}^{\text{NLOS}}\text{ATTR}_s^{\mathcal{Q}^E}$  for Top-5 beam pair selection with different location error  $\mathcal{N}(0, \sigma_c^2)$ .

the  $\text{ATTR}_s^{\mathcal{Q}^V}$  of the SABA is significantly worse than the other four methods and only has a small increase with the increase of the training epochs. This indicates that VLF cannot provide accurate beam alignment compared with the features of the other four methods.

The optimal weights of the five DNNs are determined by selecting the best epoch according to the  $\text{ATTR}_s^{\mathcal{Q}^V}$  in all epochs. Then, we compare  $\text{ATTR}_s^{\mathcal{Q}^E}$  of VBALA, VBALU, BMBA, LBA and SABA to further confirm their performance differences. Specifically, according to whether the LOS path of channel between RSU and MS is blocked by surrounding objects, we divide all the test samples into two categories: LOS samples and NLOS samples. The number of LOS and NLOS test samples are 1930 and 1410, respectively. Denote the  $\text{ATTR}_s^{\mathcal{Q}^E}$  for LOS and NLOS test samples as  ${}^{\text{LOS}}\text{ATTR}_s^{\mathcal{Q}^E}$  and  ${}^{\text{NLOS}}\text{ATTR}_s^{\mathcal{Q}^E}$ , respectively. The  ${}^{\text{LOS}}\text{ATTR}_s^{\mathcal{Q}^E}$  and  ${}^{\text{NLOS}}\text{ATTR}_s^{\mathcal{Q}^E}$  achieved by Top- $B$  beam pair selection are shown in Fig. 12 and Fig. 13, respectively. It can be seen that  ${}^{\text{LOS}}\text{ATTR}_s^{\mathcal{Q}^E}$  of VBALA is much similar to that of BMBA and LBA, but  ${}^{\text{NLOS}}\text{ATTR}_s^{\mathcal{Q}^E}$  of VBALA is better than that of BMBA and LBA. This indicates that though the hardware cost of the camera is lower than the Lidar, VBALA can outperform LBA, and VBALA can also outperform the BMBA without additional feedback overhead of MS location.

The key benefit of VBALA lies in the NLOS scenario, while the simple LOS scenario cannot cause difference between the performance of VBALA, BMBA and LBA. Compared with LBA, the advantage of VBALA under NLOS scenario can be explained as follows: VDF can contain more exact vehicle distribution information, such as the vehicle lengths, widths and heights, whereas LBA simply performs grid quantification for the scanned point cloud to design PDF. Compared with BMBA, the advantage of VBALA is not needed for MS identification. Under the NLOS scenario, since the blockage in the camera view of BS may cause the partial loss of the visual information of MS, the BMBA can not identify and distinguish well the MS and surrounding scatters from the taken image, which leads to the worse performance than the VBALA.

It is also seen that  ${}^{\text{LOS}}\text{ATTRR}_s^{\mathcal{Q}^E}$  of VBALU is slightly worse than that of VBALA. The reason is that under the LOS scenario, the optimal beam pair depends mostly on the MS's location, but the SIF of VBALU can only reflect the rough MS location by the background information of the images. However,  ${}^{\text{NLOS}}\text{ATTRR}_s^{\mathcal{Q}^E}$  of VBALU can outperform that of the VBALA, BMBA, LBA and SABA. The reason is that under the NLOS scenario, the distribution information, i.e., the location, size and orientation information, of surrounding scattering objects is more critical than the MS's location information for the beam pair selection, and there is a loss of the location information of surrounding vehicles in the VDF design owing to grid quantization.

Since SIF does not have the grid quantization loss and the VBALU does not need the MS identification, the VBALU can outperform VBALA, BMBA, LBA and SABA for NLOS scenario even without the MS's location. Nevertheless, the computational overhead of SIBAN is higher than VDBAN, since the SIBAN has about  $1.22 \times 10^{10}$  floating point operations (FLOPs) that is significantly higher than the  $3.81 \times 10^6$  FLOPs of VDBAN. The reason is that SIBAN generally requires larger scale network structure than VDBAN to process SIF whose dimensions are much larger than VDF. Moreover, the generalization of VBALU is worse than that of VBALA, as discussed in Section III.B. In addition,  ${}^{\text{LOS}}\text{ATTRR}_s^{\mathcal{Q}^E}$  and  ${}^{\text{NLOS}}\text{ATTRR}_s^{\mathcal{Q}^E}$  of SABA are the worst and are approximately 50% and 30% lower than the VBALA and VBALU, respectively, which demonstrates the design of the proposed VDF and SIF is more reasonable than VLF.

We next consider the impact of the MS location estimation error on the  ${}^{\text{LOS}}\text{ATTRR}_s^{\mathcal{Q}^E}$  and  ${}^{\text{NLOS}}\text{ATTRR}_s^{\mathcal{Q}^E}$ . We compare VBALA, VBALU, BMBA, LBA and SABA under different degrees of location estimation error to evaluate their robustness. The location error is generated as a 2D vector, and the elements of the location error are independent and set to obey the Gaussian dis-

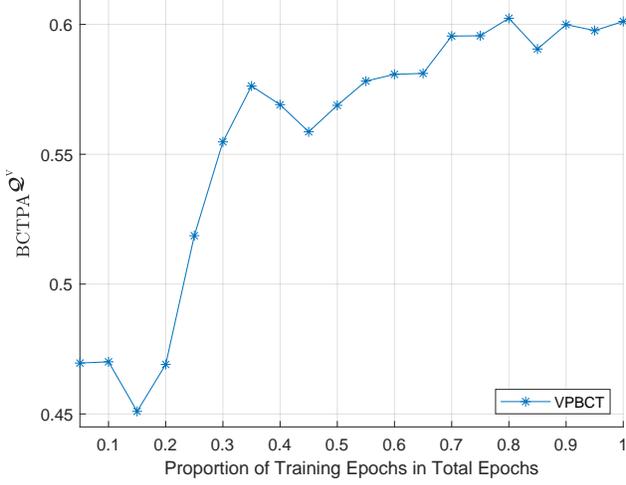


Fig. 15. BCTPA $\mathcal{Q}^v$  with the increase of the number of training epochs.

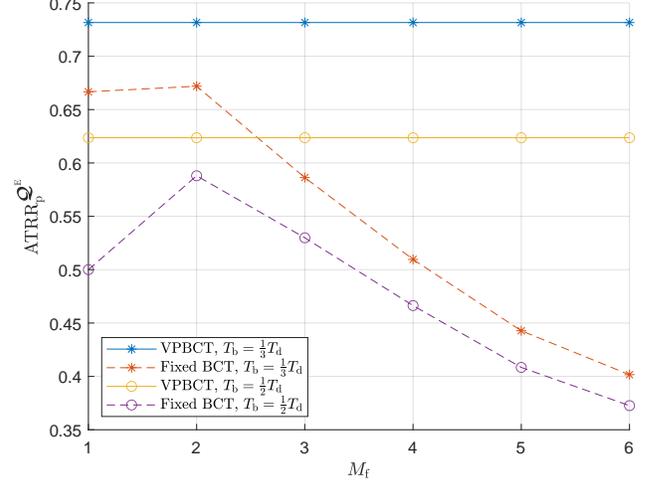


Fig. 16. ATRR $\mathcal{Q}^E$  under different  $M_f$  and  $T_b$ .

tribution  $\mathcal{N}(0, \sigma_c^2)$ . The  $\text{LOS ATRR}_s^{\mathcal{Q}^E}$  and  $\text{NLOS ATRR}_s^{\mathcal{Q}^E}$  achieved by Top-5 beam pair selection of all methods under different location error standard deviation  $\sigma_c$  are shown in Fig. 14. With the increasing of the location error variance, both  $\text{LOS ATRR}_s^{\mathcal{Q}^E}$  and  $\text{NLOS ATRR}_s^{\mathcal{Q}^E}$  of VBALA are always better than that of LBA. Hence, VBALA has better robustness than LBA under both LOS and NLOS scenarios.

It is seen that the  $\text{LOS ATRR}_s^{\mathcal{Q}^E}$  of BMBA is approximately 1.5% higher on average than the VBALA for different degrees of location estimation error. The reason may be that the images taken at BS will not be affected by the location error, whereas the VDF may change due to the location error. However, the  $\text{NLOS ATRR}_s^{\mathcal{Q}^E}$  of VBALA can be approximately 4% higher on average than the BMBA, due to the superiority of MS sensing without the need of MS identification. Thus, the VBALA's overall ATRR for both LOS and NLOS samples can be higher than BMBA, as indicated in Fig. 11. Since the difficulty of beam alignment in the NLOS scenario is significantly higher than that in the LOS scenario, as shown in Fig. 12 and Fig. 13, the LOS and NLOS detection can be pre-executed in practice to determine which method to adopt under the scenario with large location error to further improve the ATRR. Moreover, the results reveal the necessity to integrate the visual perception of BS and MS for achieving more robust beam alignment performance.

It is interesting that the performance of SABA can remain almost unchanged for LOS and NLOS scenarios, even when  $\sigma_c$  reaches 0.5m. The reason lies in that the SABA mainly adopts

$(x_M^{i,j}, y_M^{i,j}, z_M^{i,j})$ ,  $j = 1, 2, \dots, O_i$ ,  $i = 1, 2, \dots, C$ , to design VLF, and  $(x_M^{i,j}, y_M^{i,j}, z_M^{i,j})$  is obtained from 3D detection and thus is independent of the MS location error. Hence, VLF is almost unaffected by MS location error, which leads to the strong robustness of SABA. However, the performance achieved by SABA is much lower than the other methods. When  $\sigma_c > 0.16\text{m}$ , both  $\text{LOS} \text{ATTRR}_s^{\mathcal{Q}^E}$  and  $\text{NLOS} \text{ATTRR}_s^{\mathcal{Q}^E}$  of VBALA and BMBA are worse than that of VBALU. This indicates that when the location error is serious, VBALU will perform better than the other four methods.

Next, we analyze the BCT prediction performance of VPBCT. The  $\text{BCTPA}^{\mathcal{Q}^v}$  versus the increase of the number of training epochs is shown in Fig. 15. It is seen that SIBPN is trained to reach convergence and the  $\text{BCTPA}^{\mathcal{Q}^v}$  can reach about 60%. The optimal weight of SIBPN is determined according to the  $\text{BCTPA}^{\mathcal{Q}^v}$  in all epochs. We compare the  $\text{ATTRR}_p^{\mathcal{Q}^E}$  of VPBCT and the conventional fixed BCT method. The BCT adopted by the fixed BCT method is  $M_f T_d$ , where  $M_f$  is an integer with  $M_f \geq 1$ . The  $\text{ATTRR}_p^{\mathcal{Q}^E}$  of VPBCT and the fixed BCT method under different  $M_f$  and  $T_b$  are shown in Fig. 16. When  $T_b = \frac{1}{3}T_d$  or  $T_b = \frac{1}{2}T_d$ , the  $\text{ATTRR}_p^{\mathcal{Q}^E}$  of the fixed BCT method is optimal with  $M_f = 2$ , and the optimal  $\text{ATTRR}_p^{\mathcal{Q}^E}$  are 67.2% and 58.8%, respectively. However, the  $\text{ATTRR}_p^{\mathcal{Q}^E}$  achieved by the VPBCT are 73.2% and 62.4% respectively for  $T_b = \frac{1}{3}T_d$  and  $T_b = \frac{1}{2}T_d$ . This demonstrates the proposed VPBCT can get higher transmission rate than the fixed BCT method.

## VI. CONCLUSION

In this paper, we propose two beam alignment methods and a BCT prediction method with the aid of visual perception at MS. The utilization of visual ability from MS's camera can avoid the huge hardware overhead of Lidar and eliminate the privacy concerns and communication cost brought by the visual perception in BS. Under LOS scenario, the performance of VBALA, VBALU, BMBA and LBA is similar. For the NLOS scenario, VBALA and VBALU have the clear advantages in terms of hardware requirements and beam alignment accuracy compared with LBA, and can achieve better beam alignment performance without additional communication cost compared with BMBA. VPBCT can effectively improve the transmission rate by capturing the changes of communication environment from the sequence of scene images. In fact, the visual perception of both BS and MS can only obtain some local observations for the environment. Integrating with the perception ability of different devices with different sensors to enhance the communication performance deserves further research.

## REFERENCES

- [1] W. Roh et al., "Millimeter-wave beamforming as an enabling technology for 5G cellular communications: Theoretical feasibility and prototype results," *IEEE Commun. Mag.*, vol. 52, no. 2, pp. 106-113, Feb. 2014.
- [2] S. Noh, M. D. Zoltowski, and D. J. Love, "Multi-Resolution codebook and adaptive beamforming sequence design for millimeter wave beam alignment," *IEEE Trans. Wireless Commun.*, vol. 16, no. 9, pp. 5689-5701, Sept. 2017.
- [3] F. Sohrabi and W. Yu, "Hybrid digital and analog beamforming design for large-Scale antenna arrays," *IEEE J. Sel. Topics Signal Process.*, vol. 10, no. 3, pp. 501-513, Apr. 2016.
- [4] N. Gonzalez-Prelcic, A. Ali, V. Va, and R. W. Heath, "Millimeter-wave communication with out-of-band information," *IEEE Commun. Mag.*, vol. 55, no. 12, pp. 140-146, Dec. 2017.
- [5] W. Xu, F. Gao, S. Jin, and A. Alkhateeb, "3D scene based beam selection for mmWave communications," *IEEE Wireless Commun. Lett.*, vol. 9, no. 11, pp. 1850-1854, Nov. 2020.
- [6] Y. Wang, A. Klautau, M. Ribero, A. C. K. Soong, and R. W. Heath, "MmWave vehicular beam selection with situational awareness using machine learning," *IEEE Access*, vol. 7, pp. 87479-87493, 2019.
- [7] N. González-Prelcic, R. Méndez-Rial, and R. W. Heath, "Radar aided beam alignment in mmWave V2I communications supporting antenna diversity," *Proc. Inf. Theory and Appl. Workshop (ITA)*, 2016, pp. 1-7.
- [8] F. Liu, W. Yuan, C. Masouros, and J. Yuan, "Radar-assisted predictive beamforming for vehicular links: Communication served by sensing," *IEEE Trans. Wireless Commun.*, vol. 19, no. 11, pp. 7704-7719, Nov. 2020.
- [9] U. Demirhan, and A. Alkhateeb, "Radar aided 6G beam prediction: Deep learning algorithms and real-world demonstration," *arXiv preprint arXiv:2111.09676*, 2021.
- [10] M. Dias, A. Klautau, N. González-Prelcic, and R. W. Heath, "Position and LIDAR-aided mmWave beam selection using deep learning," in *Proc. IEEE Int. Workshop on Signal Processing Adv. in Wireless Commun. (SPAWC)*, Cannes, France, 2019, pp. 1-5.
- [11] A. Klautau, N. González-Prelcic, and R. W. Heath, "LIDAR data for deep learning-based mmWave beam-selection," *IEEE Wireless Commun. Lett.*, vol. 8, no. 3, pp. 909-912, 2019.
- [12] M. B. Mashhadi, M. Jankowski, T. -Y. Tung, S. Kobus, and D. Gündüz, "Federated mmWave beam selection utilizing LIDAR data," *IEEE Wireless Commun. Lett.*, vol. 10, no. 10, pp. 2269-2273, Oct. 2021.
- [13] G. Charan, M. Alrabeiah, and A. Alkhateeb, "Vision-aided 6G wireless communications: blockage prediction and proactive handoff," *IEEE Trans. Veh. Technol.*, vol. 70, no. 10, pp. 10193-10208, Oct. 2021.
- [14] Y. Tian, G. Pan, and M. -S. Alouini, "Applying deep-learning-based computer vision to wireless communications: methodologies, opportunities, and challenges," *IEEE Open J. Commun. Soc.*, vol. 2, pp. 132-143, 2021.
- [15] S. Jiang, and A. Alkhateeb, "Computer vision aided beam tracking in a real-world millimeter wave deployment," *arXiv preprint arXiv:2111.14803*, 2021.
- [16] G. Charan, T. Osman, A. Hredzak, N. Thawdar, and A. Alkhateeb, "Vision-position multi-modal beam prediction using real millimeter wave datasets," *2022 IEEE Wireless Commun. Netw. Conf. (WCNC)*, 2022, pp. 2727-2731.
- [17] V. M. De Pinho, M. L. R. De Campos, L. U. Garcia, and D. Popescu, "Vision-aided radio: User identity match in radio and video domains using machine learning," *IEEE Access*, vol. 8, pp. 209619-209629, 2020.
- [18] W. Xu, F. Gao, J. Zhang, X. Tao, and A. Alkhateeb, "Deep learning based channel covariance matrix estimation with user location and scene images," *IEEE Trans. Commun.*, vol. 69, no. 12, pp. 8145-8158, Dec. 2021.
- [19] A. M. Sayeed, "Deconstructing multiantenna fading channels," *IEEE Trans. Signal Processing*, vol. 50, no. 10, pp. 2563-2579, Oct. 2002.
- [20] <https://github.com/whxuuu/vision-communication-dataset>.

- [21] M. Daily, S. Medasani, R. Behringer, and M. Trivedi, "Self-driving cars," *Computer*, vol. 50, no. 12, pp. 18-23, December 2017.
- [22] E. Yurtsever, J. Lambert, A. Carballo and K. Takeda, "A survey of autonomous driving: Common practices and emerging technologies," *IEEE Access*, vol. 8, pp. 58443-58469, 2020.
- [23] A. Guerna, S. Bitam, and C.T. Calafate, "Roadside unit deployment in internet of vehicles systems: A survey," *Sensors*, vol. 22, no. 9, pp. 3190-3220, 2022.
- [24] S. Babu, I. Ghosh, and B. S. Manoj, "Effort: A new metric for roadside unit placement in 5G enabled vehicular networks," *2020 IEEE 3rd 5G World Forum (5GWF)*, 2020, pp. 263-268.
- [25] F. Busacca, C. Grasso, S. Palazzo, and G. Schembra, "A smart road side unit in a microeolic box to provide edge computing for vehicular applications," *IEEE Trans. Green Commun. Netw.*, 2022.
- [26] E. Arnold, O. Y. Al-Jarrah, M. Dianati, S. Fallah, D. Oxtoby, and A. Mouzakitis, "A survey on 3D Object detection methods for autonomous driving applications," *IEEE Trans. Intell. Transp. Syst.*, vol. 20, no. 10, pp. 3782-3795, Oct. 2019.
- [27] Z. Liu, Z. Wu, and R. Tóth, "Smoke: Single-stage monocular 3D object detection via keypoint estimation," in *Proc. IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2020, pp. 996-997.
- [28] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in neural information processing systems*, 2017.
- [29] J. Lu, D. Batra, D. Parikh, and S. Lee, "Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks," *arXiv preprint arXiv:1908.02265*, 2019.
- [30] W. Kim, B. Son, and I. Kim, "Vilt: Vision-and-language transformer without convolution or region supervision," *arXiv preprint arXiv:2102.03334*, 2021.
- [31] A. Prakash, K. Chitta, and A. Geiger, "Multi-Modal fusion transformer for end-to-end autonomous driving," *Proc. IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 7077-7087.
- [32] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770-778.
- [33] T. Akilan, Q. J. Wu, A. Safaei, J. Huo, and Y. Yang, "A 3D CNN-LSTM-based image-to-image foreground segmentation" *IEEE Trans. Intell. Transp. Syst.*, vol. 21, no. 3, pp. 959-971, Mar. 2020.
- [34] L. Yang and W. Zhang, "Beam tracking and optimization for UAV communications," *IEEE Trans. Wireless Commun.*, vol. 18, no. 11, pp. 5367-5379, Nov. 2019.
- [35] V. Va, X. Zhang and R. W. Heath, "Beam switching for millimeter wave communication to support high speed trains," *2015 IEEE 82nd Vehicular Technology Conference (VTC2015-Fall)*, 2015, pp. 1-5.
- [36] D. Zhang, A. Li, M. Shirvanimoghaddam, P. Cheng, Y. Li and B. Vucetic, "Codebook-based training beam sequence design for millimeter-wave tracking systems," *IEEE Trans. Wireless Commun.*, vol. 18, no. 11, pp. 5333-5349, Nov. 2019.
- [37] L. Zhang, G. Zhu, P. Shen, J. Song, S. Afaq Shah, and M. Bennamoun, "Learning spatiotemporal features using 3DCNN and convolutional LSTM for gesture recognition" *Proc. IEEE International Conference on Computer Vision (ICCV) Workshops*, 2017, pp. 3120-3128.
- [38] S. Xingjian et al., "Convolutional LSTM network: A machine learning approach for precipitation nowcasting," *Proc. Adv. Neural Inf. Process. Syst.*, 2015, pp. 802-810.
- [39] <https://carla.org>
- [40] <https://www.eclipse.org/sumo>
- [41] <https://www.remcom.com/wireless-insite-em-propagation-software>