

Environment Semantic Aided Communication: A Real World Demonstration for Beam Prediction

Shoaib Imran, Gouranga Charan, and Ahmed Alkhateeb

School of Electrical, Computer, and Energy Engineering, Arizona State University

Emails: {s.imran, gcharan, alkhateeb}@asu.edu

Abstract—Millimeter-wave (mmWave) and terahertz (THz) communication systems adopt large antenna arrays to ensure adequate receive signal power. However, adjusting the narrow beams of these antenna arrays typically incurs high beam training overhead that scales with the number of antennas. Recently proposed vision-aided beam prediction solutions, which utilize *raw RGB images* captured at the basestation to predict the optimal beams, have shown initial promising results. However, they still have a considerable computational complexity, limiting their adoption in the real world. To address these challenges, this paper focuses on developing and comparing various approaches that extract lightweight semantic information from the visual data. The results show that the proposed solutions can significantly decrease the computational requirements while achieving similar beam prediction accuracy compared to the previously proposed vision-aided solutions.

Index Terms—Millimeter wave, environment semantics, deep learning, computer vision, camera, beam selection.

I. INTRODUCTION

Future communication systems are shifting to higher frequency bands like mmWave in 5G and potentially sub-terahertz in 6G and beyond. These bands offer high bandwidth, enabling the communication systems to support the increasing data rate demands of new applications like autonomous driving, 8K video streaming, and mixed reality [1]. However, these systems require large antenna arrays and use narrow directive beams for both the transmitter and receiver to ensure sufficient receive signal power. Selecting the optimal beams for these arrays results in a large beam training overhead. This high beam training overhead makes it challenging to support highly mobile and low-latency applications, making it a key challenge in deploying these systems. Therefore, there is a need to find new ways to reduce this training overhead and enable highly-mobile mmWave/THz communication systems.

Several solutions to reduce the beam training and channel estimation overhead in mmWave/THz communication systems have been proposed over the years [2]–[4]. These solutions have included constructing adaptive beam codebooks [5], designing beam tracking techniques [3], and leveraging the channel sparsity and compressive sensing tools [2], [4]. While these classical approaches were able to achieve some improvement, they typically only save one order of magnitude in the training overhead, which is not enough for systems with large antenna arrays, particularly for serving highly-mobile and low-latency wireless communication applications.

The challenges faced by classical approaches have led to the development of machine learning-based (ML) solutions to

address the beam training and channel estimation overhead in mmWave/THz communication systems [6]–[9]. For instance, sensory data such as user position and orientation [6], RGB images [7], LiDAR point clouds [8], and radar measurements [9] can be utilized to reduce the overall training overhead. Recent work on vision-aided beam prediction has shown initial promising results [7], [10]. The current solutions use raw RGB images and CNN-based architectures to predict the optimal beams. Utilizing the raw RGB images for downstream wireless communication tasks results in higher storage requirements and increased computational cost, limiting their applicability in the real world.

Instead of directly using the raw RGB images, one promising solution can be to extract environment semantics from the images and then utilize that information to predict the optimal beam indices. Environment semantics in the form of image masks, bounding boxes, etc., present several benefits over raw RGB images. Specifically, these semantics only consist of relevant information such as object class and location, and it helps in reducing the storage cost. Further, the lower complexity of the semantics helps to reduce the computation requirements in the downstream task. An important question is whether environment semantics aided beam prediction solutions can achieve similar performance as that in [7], [10]. In this paper, we attempt to answer this important question. The main contributions can be summarized as follows:

- Formulating the environment semantic aided beam prediction problem for mmWave and THz communication systems considering practical sensing/visual and communication models.
- Developing a deep learning based solution for mmWave/THz beam prediction that utilizes different environment semantics such as image masks and bounding boxes of the detected objects
- Providing the first real-world evaluation of environment semantic-aided beam prediction based on our large-scale dataset, DeepSense 6G [11], that consists of co-existing multi-modal sensing and wireless communication data.
- Comparing the performance of various semantic-based approaches in terms of beam prediction accuracy and computational complexity, and making important conclusions on which semantic could be more useful in practice.

Based on the adopted real-world datasets, the developed solution can significantly reduce the computational complexity

while achieving similar beam prediction accuracy as compared to the prior vision-based solutions.

II. ENVIRONMENT SEMANTIC AIDED BEAM PREDICTION: SYSTEM MODEL AND PROBLEM FORMULATION

This work considers a communication system where a mmWave basestation is serving a mobile user (vehicle) in a real wireless communication environment. In this section, we first present the adopted wireless communication system model. Next, we formulate the environment semantics-aided beam prediction problem.

A. System Model

This paper adopts the system model, where a basestation, equipped with an M -element uniform linear array (ULA) and an RGB camera, is serving a mobile user. The user carries a single-antenna transmitter and is equipped with a GPS receiver capable of collecting real-time position information. The adopted communication system employs OFDM transmission with K subcarriers and a cyclic prefix of length D . To serve the mobile user, the basestation is assumed to employ a pre-defined beamforming codebook $\mathcal{F} = \{\mathbf{f}_q\}_{q=1}^Q$, where $\mathbf{f}_q \in \mathbb{C}^{M \times 1}$ and Q is the total number of beamforming vectors. Let $\mathbf{h}_k[t] \in \mathbb{C}^{M \times 1}$ denote the channel between the basestation and the mobile user at the k th subcarrier and time t , then the downlink received signal at the user can be written as

$$y_k[t] = \mathbf{h}_k^T[t] \mathbf{f}_q[t] x + v_k[t], \quad (1)$$

where $\mathbf{f} \in \mathcal{F}$ is the optimal beamforming vector at time t and $v_k[t]$ is a noise sample drawn from a complex Gaussian distribution $\mathcal{N}_{\mathbb{C}}(0, \sigma^2)$. The transmitted complex symbol $x \in \mathbb{C}$ satisfies the power constraint $\mathbb{E}[|x|^2] = P$, where P is the average symbol power. The beamforming vector $\mathbf{f}^*[t] \in \mathcal{F}$ at each time step t is selected to maximize the average receive SNR and is defined as

$$\mathbf{f}^*[t] = \underset{\mathbf{f}_q[t] \in \mathcal{F}}{\operatorname{argmax}} \frac{1}{K} \sum_{k=1}^K \text{SNR} |\mathbf{h}_k^T[t] \mathbf{f}_q[t]|^2, \quad (2)$$

where SNR is the transmit signal-to-noise ratio, $\text{SNR} = \frac{P}{\sigma^2}$.

B. Problem Formulation

Given the system model in Section II-A, at any given time instant t , the task of beam prediction can be defined as selecting the optimal beamforming vector $\mathbf{f}^*[t] \in \mathcal{F}$ that maximizes the average receive power. As presented in (2), computing the optimal beam indices require explicit channel knowledge, which is, in general, hard to acquire. One other way is to perform exhaustive search over the pre-defined beam codebook. However, the mmWave/THz communication systems need to deploy large antenna arrays and use narrow directed beams to guarantee sufficient receive SNR. Selecting the optimal beams for these systems with large antenna arrays through exhaustive search is typically associated with large training overhead; making it challenging for these systems to support high mobility wireless communication applications.

One promising way to reduce the large beam training overhead is to develop machine learning-based solution that leverages prior observation and additional side information for fast mmWave/THz beam prediction. Recent work on sensing-aided beam prediction have achieved initial promising results by utilizing sensory data such as GPS positions [6], RGB images [7], LiDAR point clouds [8] and radar observations [9]. In this paper, we propose to utilize additional sensory data (RGB images captured by the basestation) to predict the optimal index for the transmitter in the scene. However, instead of using the raw images captured by the basestation, we propose to extract environment semantics (object masks, bounding boxes, etc.) from the RGB images. These extracted semantics can then be used to predict the optimal beam indices.

In this work, we target predicting the optimal beam indices based on the availability of RGB images captured by a camera installed at the basestation. Formally, we define $\mathbf{X}[t] \in \mathbb{R}^{W \times H \times C}$ as the corresponding RGB image, captured by a camera installed in the basestation at time t , where W , H , and C are the width, height, and the number of color channels of the image. Let, $\mathbf{S}[t]$ represent the environment semantics extracted from the visual data captured by the mmWave/THz basestation. The objective of the beam prediction task is to find a mapping function f_{Θ} that utilizes the available semantics, $\mathbf{S}[t]$ to predict the (estimate) optimal beam index $\hat{\mathbf{f}}[t] \in \mathcal{F}$ with high fidelity. The mapping function can be formally expressed as

$$f_{\Theta} : \mathbf{S}[t] \rightarrow \hat{\mathbf{f}}[t]. \quad (3)$$

In this work, we develop a machine learning model to learn this prediction function f_{Θ} . Let $\mathcal{D} = \{(\mathbf{S}_u, \mathbf{f}_u^*)\}_{u=1}^U$ represent the available dataset consisting of image-beam pairs is collected from the real wireless environment, where U is the total number of samples in the dataset. Then, the goal is to maximize the number of correct predictions over all the samples in the dataset \mathcal{D} . This can be formally written as

$$f_{\Theta}^* = \underset{f_{\Theta}}{\operatorname{argmax}} \prod_{u=1}^U \mathbb{P}(\hat{\mathbf{f}}_u = \mathbf{f}_u^* | \mathbf{S}_u). \quad (4)$$

The prediction function is parameterized by a set of model parameters Θ and is learned from the labeled data samples in the dataset \mathcal{D} . The objective is to find the best parameters Θ^* that maximize the product of the probabilities of correct predictions. Next, we present our proposed machine learning-based solution for environment semantics-aided mmWave/THz beam prediction.

III. ENVIRONMENT SEMANTIC AIDED BEAM PREDICTION: THE KEY IDEA AND A DEEP LEARNING SOLUTION

This section presents an in-depth overview of the proposed environment semantics-aided beam prediction solution. For this, we first present the key idea in Section III-A followed by the details of our proposed solution in Section III-B.

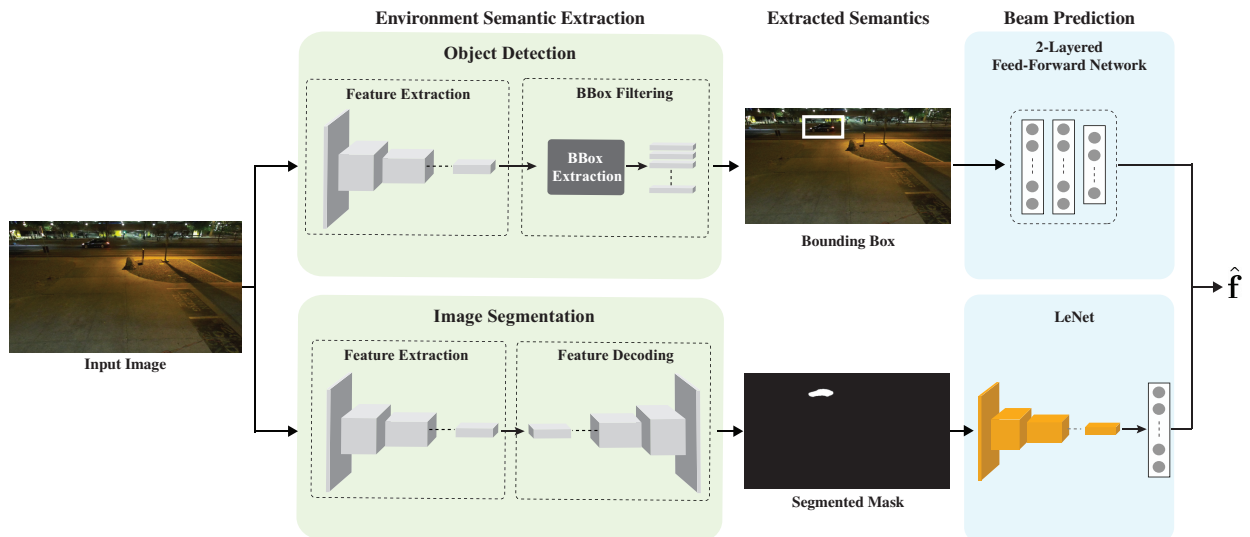


Fig. 1. A block diagram showing the proposed solution for environment semantics-aided beam prediction task. As shown in the figure, the camera installed at the basestation captures real-time images of the wireless environment. We propose to first extract environment semantics (object masks, bounding boxes, etc.) from the RGB images. These extracted semantics can then be used to predict the optimal beam indices.

A. The Key Idea

Enabling highly-mobile mmWave/THz wireless communication applications requires overcoming some of the critical challenges associated with these high-frequency systems. One key challenge arises from the severe path loss associated with these high-frequency signals. In order to overcome this challenge, the mmWave/THz communication systems adopt large antenna arrays and use narrow directive beams in both the transmitter and receiver. However, selecting the optimal beams in these systems with large antenna arrays results in large beam training overhead, making it challenging to support high-mobility wireless communication applications. In general, the mechanism of directing the narrow beams can be viewed as focusing the wireless signal in a particular direction in space. The beam vectors theoretically divide the wireless environment (spatial dimension) into multiple (possibly overlapping) sectors. Therefore, if we have access to a pre-defined codebook, the beam prediction task can be re-defined as a classification task, where depending on where the user is located in the wireless environment, a particular beam index from the beam codebook can be assigned.

In order to perform the beam classification task, it is imperative to understand and extract the exact user's location in the wireless environment. One promising way to achieve this is by utilizing additional sensing modalities such as vision, GPS, etc. The recent advancements in machine learning and computer vision, in particular, have provided capabilities to accurately identify and locate the objects of interest in the visual scene. For example, object detection models such as You Only Look Once version 7 (YOLOv7) can detect different objects in the visual data and provide bounding box coordinates. This work utilizes visual data to reduce the beam training overhead associated with mmWave/THz communication systems. Prior work on vision-aided beam prediction for

mmWave/THz communication systems [7], [10] has primarily proposed solutions that provide the image as input directly to a convolutional neural network to predict the optimal beams. Although these initial approaches have provided promising results, the proposed solutions are computationally expensive, which further limits their adoption in the real world. In this work, we propose to first extract environment semantics from the images (captured at the basestation), such as image segmentation mask and bounding box information. The extracted semantics is then utilized to predict the optimal beamforming vector from a pre-defined beam codebook.

B. Proposed Solution

In this subsection, we propose a two stage environment semantics aided beam prediction solution as shown in Fig. 1. In the first stage, a machine learning model is deployed to extract the environment semantics in the form of either masks or bounding boxes of the objects of interest. The semantics are then utilized, in the second stage, to predict the beam index from a pre-defined beamforming codebook.

Environment Semantics Extraction: In the first stage, we adopt a deep neural network (DNN) that takes the RGB image \mathbf{X} as input and provides the semantic representation for the RGB image, \mathbf{S} , as the output. In general, the DNN must fulfill two essential conditions: (i) It should provide an accurate semantic representation and (ii) have a low computational footprint. In this paper, we adopt two state-of-the-art object detection models, YOLOv7 [12], and MobileNet version 2 (MobileNetv2) [13] for generating these semantics. It is important to note that the YOLOv7 model has significantly more parameters than MobileNetv2. In addition, while YOLOv7 and MobileNetv2 have been designed specifically for object detection and generating bounding box coordinates of the detected objects, they can be adapted for image segmentation

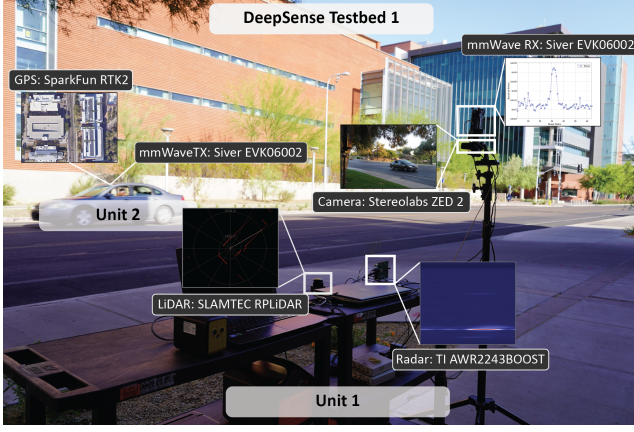


Fig. 2. This figure presents the DeepSense 6G testbed 1 and the different sensing modalities. It consists of two units: Unit 1 (a stationary unit), acting as the basestation, and unit 2 (a vehicle), representing a mobile user.

as demonstrated in [14]. To that end, we design a semantic segmentation head on top of a MobileNetv2 architecture for extracting the masks of the detected objects. Moreover, instead of training these models from scratch, we utilize pre-trained models as they can detect most of the relevant objects found in a wireless environment. To perform a detailed comparison between accuracy and latency, we generate two different types of semantics to represent the users in the environment, namely their binary masks $\mathbf{S}_{\text{mask}} \in \mathbb{R}^{W \times H}$ and their bounding box vectors $\mathbf{S}_{\text{bbox}} \in \mathbb{R}^{4 \times 1}$. The bounding box vector $\mathbf{S}_{\text{bbox}} = [x_c, y_c, w, h]$ where x_c, y_c, w , and h are the x -centre, y -center, width, and height of the detected object, respectively. The extracted binary mask is then downsampled to the size $\hat{W} \times \hat{H}$ before it is provided as input to the DNN in the second stage. Unlike RGB images, the binary masks mostly retain their meaningful information, namely the positions and the shapes of the detected objects, when they are downsampled. Further, both the binary masks and bounding box vectors are normalized to the range $[0, 1]$.

Beam Prediction: The second stage consists of a separate decision network that takes the semantic representation \mathbf{S} as the input and predicts the optimal beam index \hat{f} as the output. Similar to the requirements of the semantics extraction network, the decision network is expected to (i) accurately predict the beam index and (ii) have a low computational footprint. We further observe that the structure of the semantic representation is significantly different for masks and bounding box vectors. Consequently, we adopt two different DNN architectures for the two modes of semantic representation utilized in this work.

1) *Fully Connected Neural Network for \mathbf{S}_{bbox} :* For the bounding box vectors, we use a 2-layered fully connected neural network (FCNN) with 175 neurons in each layer. FCNNs work well with structured data as they use the network's weights to model the relationship between every input element. Moreover, FCNNs make dense connections between the neurons of adjacent layers enabling them to learn more complex relationships between the input elements.

2) *LeNet for \mathbf{S}_{mask} :* For images, however, convolutional neural networks (CNNs) have achieved better performance and robustness as they take advantage of the spatial correlation between the neighboring pixels. Therefore, for the binary masks, we adopt a simple CNN model (LeNet [15]) to predict the optimal beam index. The LeNet consists of 2 convolutional layers followed by two fully connected layers. The LeNet and the FCNN take masks and bounding vectors as inputs, respectively, and learns to predict the optimal beam indices.

IV. TESTBED DESCRIPTION AND DEVELOPMENT DATASET

To test the effectiveness of the proposed environment semantic-aided beam prediction solution, we utilize DeepSense 6G [11] dataset. DeepSense 6G is a real-world multi-modal dataset developed for sensing-aided wireless communication applications. It consists of co-existing multi-modal data such as mmWave wireless communication, GPS data, vision, Radar, and LiDAR collected in a real-world wireless environment. In this section, we first present a brief overview of the DeepSense 6G testbed utilized during the data collection. Next, we present an analysis of the final development dataset used for developing and evaluating the proposed beam prediction solution.

DeepSense 6G: [Scenarios 5 and 7] In this work, we adopt scenarios 5 and 7 of the DeepSense 6G dataset to evaluate the efficacy of the proposed solution. The hardware testbed and the example image samples from both scenarios are presented in Fig. 3. The DeepSense testbed 1 is utilized for this data collection, which consists of a stationary and mobile unit. The mobile unit (vehicle), acting as the transmitter, is equipped with a 60GHz quasi-omni antenna and a GPS receiver to record the real-time location of the user. The stationary unit (basestation) is equipped with (i) an RGB camera and (ii) a 16-element 60GHz mmWave phased array. It uses an over-sampled predefined codebook of 64 beams for receiving the transmitted signal. The data collected at each time instant consists of the GPS position of the user, RGB images, and the mmWave receive power vector. For further details, the reader is referred to [11], which describes the data collection testbed in detail.

DeepSense 6G Development Dataset: We separately test the effectiveness of the proposed solution on scenarios 5 and 7 of the DeepSense 6G dataset. The measurements of scenario 5 and scenario 7 are collected at different times of the day in different locations as shown in Fig. 3. In particular, scenario 7 measurements are collected in daylight whereas scenario 5 measurements are taken at night. In addition, both scenario 5 and scenario 7 contain many such images in which there are multiple mobile vehicle units but only one of them is transmitting to the basestation. Scenario 5 and scenario 7 contain 2300 and 854 samples respectively, which are further split into training, validation and testing samples with a ratio of 70/20/10 respectively.



Fig. 3. This figure presents the overview of the different data collection locations. Fig. (a) and (b) present the visual data (RGB images) captured in scenario 5 and 7, respectively. We also show the corresponding bounding box of the mobile unit in the image. Fig. (c) shows the corresponding mask of the mobile unit in the image.

TABLE I
BEAM PREDICTION: DESIGN AND TRAINING HYPER-PARAMETERS

Parameters	Mask	Bounding Box
ML Model	LeNet-5	2-layered MLP
Batch Size	64	128
Learning Rate	1×10^{-3}	1×10^{-2}
Learning Rate Decay	epochs 10 and 20	epochs 15 and 30
LR Reduction Factor	0.1	0.1
Total Training Epochs	30	50

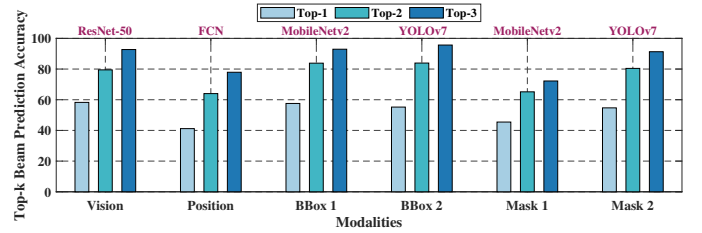
V. PERFORMANCE EVALUATION

In this section, we study the performance of the proposed solution for the environment semantic-aided beam prediction task. For this, in Section V-A, we present the details of the experimental setup adopted in this work. Next, we discuss the performance of the proposed solution in Section V-B.

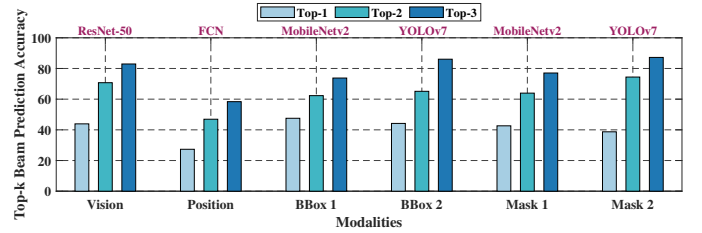
A. Experimental Setup:

In this work, we first extract the environment semantics which are then utilized to predict the optimal beam indices. As presented in Section III, we utilize two different state-of-the-art object detection models (YOLOv7 and MobileNetv2) to extract the semantics from the visual data captured by the basestation. The two types of semantics (image mask and bounding boxes) extracted from the images require specialized machine learning-based models to perform the beam prediction task. For the binary mask-based approach, we design a CNN-based architecture (similar to that of LeNet). For the bounding box-based solution, we adopt a 2-layered fully-connected neural network. The proposed ML models are trained and validated on the task-specific dataset as presented in Section IV. The cross-entropy loss with the Adam optimizer is used to train the models. The details of the hyper-parameters used to fine-tune the models are presented in Table I.

Evaluation Metric: The primary metric adopted to evaluate the proposed solution is the top- k accuracy. Note that the top- k accuracy is defined as the percentage of the test samples where the optimal ground-truth beam is within the top- k predicted beams. This work presents the top-1, top-2, and top-3 accuracies to evaluate the proposed solutions comprehensively.



(a) DeepSense Scenario 5



(b) DeepSense Scenario 7

Fig. 4. This figure plots the top- k accuracies ($k \in (1, 2, 3)$) for the proposed environment semantics-aided beam prediction solution. We, further, plot the beam prediction accuracies for two prior work with visual and position data. It is observed the proposed bounding box-based solution (MobileNetv2) achieves similar or better performance than the prior vision-based solution.

B. Numerical Results

With the experimental setup described in Section V-A, in this subsection, we study the beam prediction performance of the proposed solution.

Can environment semantics be utilized to predict the optimal beams? In order to perform a comparative study, as presented in Section III, we utilize two state-of-the-art object detection models (YOLOv7 and MobileNetv2) to extract the environment semantics. To further facilitate a detailed study, we extract the image segmentation masks and the bounding boxes of the detected objects. In Fig. 4, we present the top-1, top-2, and top-3 accuracies achieved by the different models utilized and the different extracted semantics for both scenarios 5 and 7 dataset. We also present the beam prediction accuracies of two prior work that utilize additional sensing data such as vision [7] and position [6]. It is important to highlight here that different from our approach, the prior work with

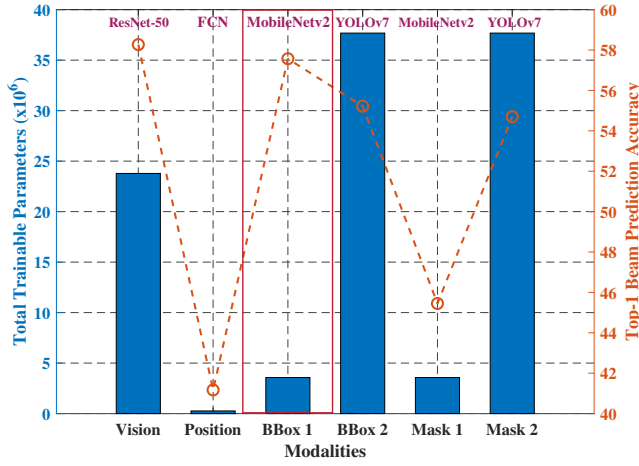


Fig. 5. This figure plots the total number of trainable parameters versus the top-1 beam prediction performance for the different approaches. It shows that the bounding box-based solution can achieve similar top-1 beam prediction accuracy while using less than $1/6^{th}$ of the parameters when compared with the baseline vision-based approach.

visual data provides the image directly as an input to a CNN. The beam prediction performance of these approaches will act as the baseline for our proposed solution. From Fig. 4, it is observed that the bounding box-based solution can achieve similar beam prediction accuracy or even surpass the vision-based baseline solution. We observe similar performance trend for both the bounding boxes extracted from YOLOv7 and MobileNetv2. The proposed solution based on image segmentation mask achieves slightly lower performance as compared to the baseline vision and bounding box-based approach. This can be attributed to the fact that the generated masks were noisy in nature; a drawback that can be improved with further processing. However, both the mask and bounding box-based solutions were able to surpass the position-alone solution, highlighting the efficacy of utilizing environment semantics for predicting the optimal beam indices in mmWave/THz communication systems.

Trade-off between computational cost and model performance: The final adoption of any proposed solution depends on achieving both high accuracy and low latency. As shown in Fig. 4, the proposed environment semantics-aided beam prediction (bounding box-based) solution achieves high prediction accuracy compared to the baseline vision-based solution. In Fig. 5, we plot the computational requirement versus the top-1 beam prediction performance for the different approaches. Although the position-based beam prediction approach has the lowest footprint (in terms of model parameters), it achieves the lowest beam prediction performance. This can be attributed to the inherent errors in GPS data, and access to accurate positions in the future might help improve the performance. It is also interesting to note that the bounding box-based approach (extracted using MobileNetv2) achieves similar performance as the baseline vision-based solution with only a fraction of the trainable model parameters. These results highlight the computational efficacy of the proposed environment semantics-based

beam prediction solution. It also emphasizes that there is an inherent trade-off between inference accuracy and achievable latency, which must be considered during network design.

VI. CONCLUSION

This paper develops a two stage deep learning based solution for fast and accurate mmWave/THz beam prediction. In the first stage, the semantics are extracted from the visual data. These semantics are then used for beam prediction in the second stage. We evaluate the proposed solution on the DeepSense 6G dataset. Using bounding box vectors from MobileNetv2, the proposed solution achieves similar top-1 beam prediction accuracy while using less than $1/6$ of the parameters when compared with the vision based approach, highlighting a promising solution for real-world systems.

REFERENCES

- [1] T. S. Rappaport, Y. Xing, O. Kanhere, S. Ju, A. Madanayake, S. Mandal, A. Alkhateeb, and G. C. Trichopoulos, "Wireless communications and applications above 100 GHz: Opportunities and challenges for 6G and beyond," *IEEE Access*, vol. 7, pp. 78 729–78 757, 2019.
- [2] A. Alkhateeb, O. El Ayach, G. Leus, and R. Heath, "Channel estimation and hybrid precoding for millimeter wave cellular systems," *IEEE Journal of Selected Topics in Signal Processing*, vol. 8, no. 5, pp. 831–846, Oct. 2014.
- [3] S. Jayaprakasam, X. Ma, J. W. Choi, and S. Kim, "Robust beam-tracking for mmwave mobile communications," *IEEE Communications Letters*, vol. 21, no. 12, pp. 2654–2657, 2017.
- [4] R. W. Heath, N. Gonzalez-Precic, S. Rangan, W. Roh, and A. M. Sayeed, "An overview of signal processing techniques for millimeter wave mimo systems," *IEEE journal of selected topics in signal processing*, vol. 10, no. 3, pp. 436–453, 2016.
- [5] Y. Zhang and A. Alkhateeb, "Learning reflection beamforming codebooks for arbitrary RIS and non-stationary channels," *arXiv preprint arXiv:2109.14909*, 2021.
- [6] J. Morais, A. Behboodi, H. Pezeshki, and A. Alkhateeb, "Position aided beam prediction in the real world: How useful GPS locations actually are?" 2022. [Online]. Available: <https://arxiv.org/abs/2205.09054>
- [7] G. Charan, T. Osman, A. Hredzak, N. Thawdar, and A. Alkhateeb, "Vision-position multi-modal beam prediction using real millimeter wave datasets," in *2022 IEEE Wireless Communications and Networking Conference (WCNC)*, 2022, pp. 2727–2731.
- [8] S. Jiang, G. Charan, and A. Alkhateeb, "LiDAR aided future beam prediction in real-world millimeter wave v2i communications," *IEEE Wireless Communications Letters*, Oct. 2022. [Online]. Available: <https://arxiv.org/abs/2203.05548>
- [9] U. Demirhan and A. Alkhateeb, "Radar aided 6G beam prediction: Deep learning algorithms and real-world demonstration," 2021. [Online]. Available: <https://arxiv.org/abs/2111.09676>
- [10] G. Charan and A. Alkhateeb, "User identification: The key enabler for multi-user vision-aided wireless communications," 2022. [Online]. Available: <https://arxiv.org/abs/2210.15652>
- [11] A. Alkhateeb, G. Charan, T. Osman, A. Hredzak, J. Morais, U. Demirhan, and N. Srinivas, "Deepsense 6G: A large-scale real-world multi-modal sensing and communication dataset," 2022. [Online]. Available: <https://arxiv.org/abs/2211.09769>
- [12] C.-Y. Wang, A. Bochkovskiy, and H.-Y. M. Liao, "Yolov7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors," *arXiv preprint arXiv:2207.02696*, 2022.
- [13] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "Mobilenetv2: Inverted residuals and linear bottlenecks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 4510–4520.
- [14] E. Mohamed, A. Shaker, A. El-Sallab, and M. Hadhoud, "Insta-yolo: Real-time instance segmentation," *arXiv preprint arXiv:2102.06777*, 2021.
- [15] "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.