

UnLoc: A Universal Localization Method for Autonomous Vehicles using LiDAR, Radar and/or Camera Input

Muhammad Ibrahim¹, Naveed Akhtar¹, Saeed Anwar², and Ajmal Mian¹

Abstract—Localization is a fundamental task in robotics for autonomous navigation. Existing localization methods rely on a single input data modality or train several computational models to process different modalities. This leads to stringent computational requirements and sub-optimal results that fail to capitalize on the complementary information in other data streams. This paper proposes UnLoc, a novel unified neural modeling approach for localization with multi-sensor input in all weather conditions. Our multi-stream network can handle LiDAR, Camera and RADAR inputs for localization on demand, i.e., it can work with one or more input sensors, making it robust to sensor failure. UnLoc uses 3D sparse convolutions and cylindrical partitioning of the space to process LiDAR frames and implements ResNet blocks with a slot attention-based feature filtering module for the Radar and image modalities. We introduce a unique learnable modality encoding scheme to distinguish between the input sensor data. Our method is extensively evaluated on Oxford Radar RobotCar, ApolloSouthBay and Perth-WA datasets. The results ascertain the efficacy of our technique.

I. INTRODUCTION

Vehicle localization in outdoor environment is an essential task in robotics, especially in the autonomous driving domain. To achieve self-autonomy in urban outdoor environment, a vehicle must be able to precisely localize itself. Current outdoor localization systems rely on the Global Navigation Satellite System (GNSS). However, the lack of accuracy and signal blockage in densely populated regions for GNSS make it an inadequate technology for autonomous vehicles. Creating an offline map of the environment and using query frames during online navigation provides a viable alternate solution to the problem. Conventional methods in this direction [1], [2] employ frame registration for localization. However, this leads to impractical computational requirements. More recently, deep learning techniques have shown great promise in addressing the issue [3].

Among the deep learning methods, 3D point cloud regression-based approaches, e.g., [3], [4], can directly predict six degrees of freedom (6DoF) poses to localize vehicles. Point clouds provide depth information of the scene and 360°

field of view (FoV), which are helpful for precise localization. However, LiDAR data is also inherently prone to high level of noise in rainy and foggy weather. Moreover, its unstructured nature demands complex and computationally expensive modeling when outdoor localization completely relies on the LiDAR data.

In the related literature, processing RGB camera images with deep learning is also considered suitable for localization. Currently, techniques such as PoseNet and its variants [5], [6] use a single or a series of image frames to predict 6DoF poses. Whereas image modality offers detailed spatial information, it is easily influenced by environmental variations, such as sunlight, rain, fog etc., which is detrimental for localization. Comparatively, Radar data is largely insensitive to various weather conditions, e.g., darkness, fog, snow and sunlight. Leveraging that, Cen *et al.* [7] extracted features from Radar scans and then applied scan matching to predict ego-motion. Radarloc [8] is a recent deep learning localization method that predicts global poses from Radar data. Nevertheless, Radar data does not have precise 3D information and is noisy, which compromises the overall localization performance.

For the applications like self-driving vehicles, robust outdoor localization is only possible by leveraging complementary characteristics of different data modalities. In this work, we present a multi-sensor localization approach, shown in Fig. 1, that learns a unified neural model, called UnLoc, for point cloud, Radar and image data, to achieve precise 6DoF localization. The proposed model uses sparse 3D convolutions to process cylindrical representations of point clouds, while 2D convolutions and slot attention-based feature filtering are used to process Radar and image modalities. We also introduce a learnable modality encoding technique to optimally discriminate between different data modalities for their on demand use. Our method allows the use of a single or any combination of modalities during inference. This makes our method robust to sensor failure.

Our network is trained on six sensor inputs for three different modalities. Due to the unique universal nature of our method, we also devise a technique to generate common ground-truth for different sensor data streams, which enables effective training of our model. Our method is an adaptive deep-learning technique that can be used in challenging real-world environments. We establish the performance of our approach on three publicly available datasets: Oxford Radar RobotCar [9], Appolo-Southbay [10], and Perth-WA [3]. We also conduct a thorough ablation study to analyze the effects of using various sensor inputs with our model. We establish

Professor Ajmal Mian is the recipient of an Australian Research Council Future Fellowship Award (project number FT210100268) funded by the Australian Government. Dr. Naveed Akhtar is recipient of an Office of National Intelligence National Intelligence Postdoctoral Grant (project number NIPG-2021-001) funded by the Australian Government.

¹ Department of Computer Science, The University of Western Australia. muhammad.ibrahim@research, naveed.akhtar@, ajmal.mian@) uwa.edu.au

²Information and Computer Science, King Fahad University of Petroleum and Minerals (KFUPM), Dhahran, KSA, saeed.anwar@kfupm.edu.sa

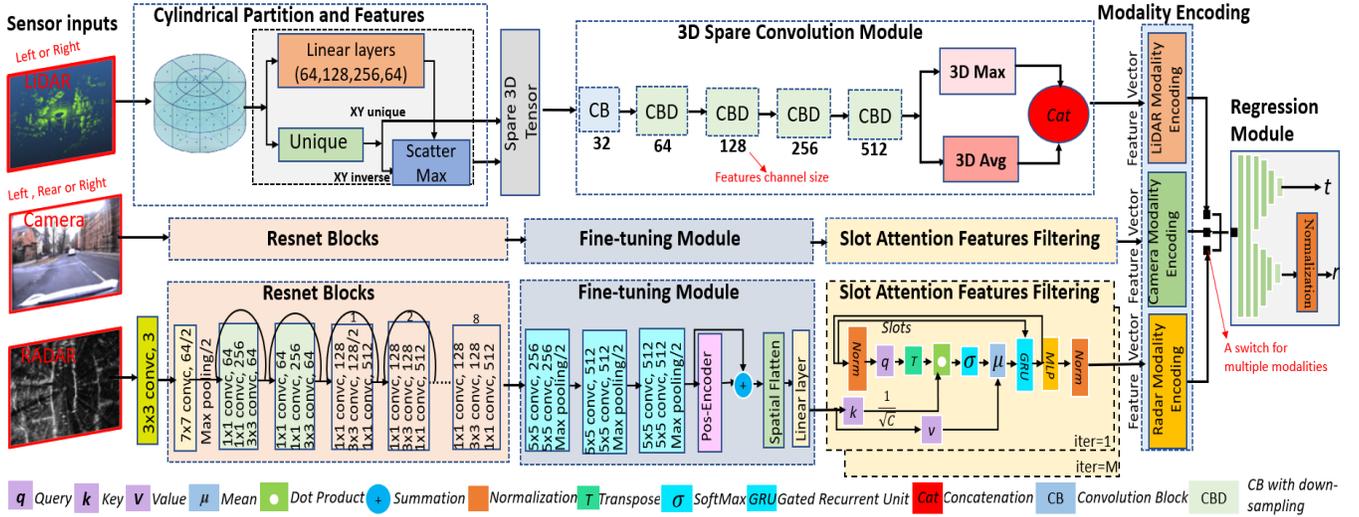


Fig. 1. Architecture of the proposed UnLoc. **Top** stream of layers of our network processes 3D point cloud data. It transforms the input into a cylindrical representation which is processed by sparse 3D convolution blocks, CB and CBD (see Fig. 2), to extract point cloud features. These features are passed through 3DMax and 3DAvg layers to compute the feature vectors. **Middle** stream of the network processes images to extract features with ResNet blocks, which are further processed by our fine tuning and a transformer based features filtering module. **Bottom** stream processes Radar modality and has the same architecture as the middle stream, except with an additional 2D Conv layer. Feature vectors for each modality are encoded by our learnable modality encoding scheme which passes the features to a Regression Module for 6DoF poses prediction.

strong localization results on all the datasets, outperforming the existing methods on each modality and further improving the performance by processing multi-modality input.

II. RELATED WORK

Localization is essential for autonomous driving [11], [12]. In the existing literature, different input data modalities, such as point cloud, image, and Radar data, distinguish different localization methods. Conventional point cloud matching techniques employ registration methods [1], [13]. On the other hand, recent point cloud deep learning-based techniques associate input frames to a point cloud map for 6DoF pose estimation [4], [10], [14], [15]. Among the mentioned methods, PointLoc [4], Slice [3] and L3Net [10] compress the map into a neural 6DoF pose predictor for vehicles. The PointLoc [4] exploits PointNet framework to predict poses while Slice [3] uses transformer architecture [16]. In general, directly predicting 6DoF pose from a point cloud frame is challenging because the unstructured LiDAR representation conflicts with the high precision demands of the task. Therefore, other methods, e.g., [5], [6] often augment their neural models to handle the data complexity.

Some works also devise camera-based deep learning localization methods. For instance, PoseNet [5] is the pioneering technique that utilizes camera images to predict 6DoF poses. Likewise, the recent variants [6], [17] of PoseNet regress poses using a single or multiple images, exploiting geometric loss and modeling poses with Bayesian Neural Network (BNN) to enhance the performance. Walch *et al.* [18] employed LSTMs for matching geometric features to improve the pose precision. Retrieval-based learning approaches, e.g., CamNet [19], RelocNet [20], and Camera Relocalization CNN [21], use agents that have previously visited the exact location. However, the above approaches are restricted in performance due to the demerits of the visual sensors.

Contemporary localization strategies also explore the Radar data modality. For example, Barnes *et al.* [22] proposed a deep correlative scan matching technique based on learned feature embedding and a self-supervised module for Radar odometry system. Later, the authors developed a deep key point detector and metric localizer [23] for Radar odometry estimation. Cen *et al.* [7] extracted features from Radar scans and then applied scan matching to predict ego motion. RadarLoc [8] is the latest method that predicts global absolute poses with respect to the world coordinate system employing Radar data. Compared to camera and LiDAR, radar is not as sensitive to the weather conditions and provides 360° FoV of a scene. However, it lacks precise 3D information when compared to the LiDAR data.

The methods mentioned above mainly rely on a single sensor data modality. However, outdoor localization, especially for autonomous driving, requires the precision and robustness that is hard to achieve with a single modality. Still, localization through multiple sensors (Radar, camera and LiDAR) in outdoor environments is currently a largely unexplored direction. In this work, we fill this gap by devising a universal localization neural model that leverages point cloud, Radar and/or camera input on demand, to provide robust and precise 6DoF pose predictions.

III. METHODOLOGY

We propose a multi-sensors localization method that can leverage point cloud, Radar and/or image modalities on demand. Our approach uses three parallel streams for these modalities in the early stages of the model, see Fig. 1. It applies sparse 3D convolutions on the LiDAR data and employs ResNet blocks followed by a slot attention-based feature filtering for the Radar and image modalities. Moreover, it employs a learnable modality encoding that learns to identify the input sensor data stream for optimal performance. The

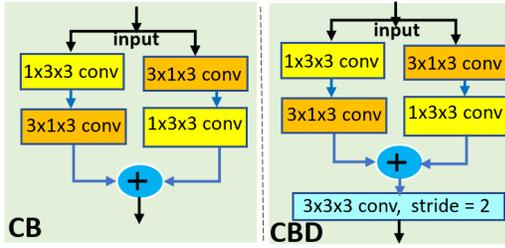


Fig. 2. Sparse 3D convolution block without down-sampling (CB) and with down-sampling (CBD, using an extra 3D convolution with stride 2).

architecture for the Radar and image data streams are largely similar, except for an additional 2D convolution layer to process the Radar data. Our technique computes a feature vector for each data modality and applies modality encoding to that. The individual data streams get projected onto their respective feature spaces that have similar dimensionality across modalities. The data features eventually get processed by a regression module to predict the 6DoF pose. Our model is trained on six sensor inputs: LiDAR (left and right), camera (left, rear and right), and Radar. For inference, it accepts any single input or any combination of the input modalities. Our method is designed for localization in outdoor environment, particularly for self-driving vehicles. Due to its multi-modality nature, it is well-suited to different weather conditions and is robust to sensor failures. Each component of our framework is explained in detail below.

A. 3D Features Extraction

We propose a sparse 3D convolution localization block to process the point cloud data, see Fig. 1. This block extracts 3D geometric features from LiDAR data, which are particularly relevant for localization in the outdoor environment. The 3D convolution using voxelization is known for its efficacy to process LiDAR data [24]. However, voxel processing using 3D convolution is computationally expensive. Hence, we devise a lightweight sparse 3D convolutional method utilizing cylindrical partitioning. Our approach uses a series of sparse 3D convolutional sub-blocks for processing point cloud data, followed by Max-pooling and Average pooling layers. Details of this process are provided below.

1) *3D Cylindrical Partition and Features*: Outdoor LiDAR point cloud has the attribute of changing density, with nearby regions having substantially higher point densities than the far regions. Hence, cylindrical coordinate system is well-suited to partition LiDAR data, which evenly distributes the points across various partitions by providing larger volumes for the far-off points. The top-left corner in Fig. 1 illustrates the workflow of our partitioning, followed by point feature processing. First, the (X, Y, Z) coordinates of points are converted to (R, θ, z) for a cylindrical grid representation, where R, θ are the radius and height, respectively. Then, a cylindrical partitioning is performed to generate voxels having largely uniform point distribution. The farther regions have larger voxels compared to the nearby regions. Next, the cylindrical grid representation is passed through an MLP-based module with linear layers to obtain cylinder

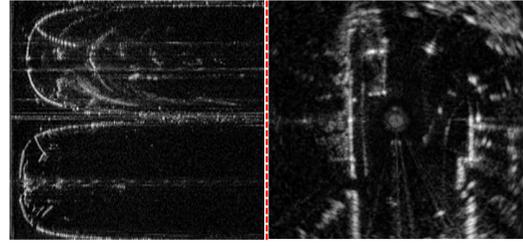


Fig. 3. Radar sensor output in polar form (Left) and after transformation to Cartesian coordinates (Right).

features. These features are further processed using Unique indices and Scatter-Max layers to obtain maximum magnitude features. Finally, we get the 3D cylinder representations with size $\in \mathbb{R}^{D \times H \times W \times L}$, where D is the feature dimension, and H, W, L respectively represent the height, width, and length of the cylinder. Our 3D Sparse convolution module subsequently processes this representation.

2) *3D Sparse Convolution*: To process the cylindrical representation, we had two options: conventional 3D convolution or sparse 3D convolution. We choose the latter for memory and computational efficiency. Inspired by Cylinder3D [24], we create two asymmetric types of 3D sparse convolution modules, without down-sampling (CB) and with down-sampling (CBD), as shown in Fig. 2. Both convolution blocks have two sparse convolution streams with stride 1. The first stream has kernel $(1, 3, 3)$ followed by kernel $(3, 1, 3)$ and the second stream has the same kernels in the reverse order. The output of both streams is added. In the convolution with down-sampling (CBD), an additional 3D convolution with kernel $(3, 3, 3)$ and stride 2 is applied for downsampling. We employ one CB and four CBD modules in series with 32, 64, 128, 256, 512 output channels. This block's output features size is $B, 512, 30, 23, 8$, representing batch size B , the feature dimensionality, and cylinder height, width and length, respectively.

3) *3D Max-pooling and Average Pooling*: To aggregate the information, we use max- and average-pooling techniques. The pooling layers compute maximum and average feature values along the spatial/cylindrical dimensions. These layers generate outputs of size $\mathbb{R}^{B \times 512}$, which are concatenated to generate the final feature vector of size $\mathbb{R}^{B \times 1024}$ for further processing, which is discussed in Sec. III-C.

B. 2D Features Extraction

Here, we explain the 2D feature extraction blocks for image and Radar modalities. The architecture for both blocks in our framework is largely identical, except for an additional 2D convolution layer used for the Radar modality, see Fig. 1. This additional layer broadcasts channel size from one to three for later processing. We also transform the polar scanning Radar outputs into Cartesian coordinates as grey-scale images $\in \mathbb{R}^{H \times W}$, as shown in Fig. 3. This transformation helps in localization performance. The architecture for our feature extraction blocks includes a ResNet module, a fine-tuning module and a slot attention-based features filtering module. These components are discussed below in detail.

1) *ResNet Blocks*: The primary responsibility of this module is to extract useful local features from Radar and camera modalities. The state-of-the-art camera-based localization methods [25], [17], [26] utilize a pre-trained ResNet model [27] as a features extractor. The aforementioned approaches typically involve selecting layers from a significant portion of the pre-trained model, resulting in computationally demanding models. Our framework focuses solely on the ten initial blocks of the pre-trained ResNet-152 model, thereby considerably reducing the computational cost. The input to this module is in $\mathbb{R}^{B \times D \times H \times W}$, where B is the batch size, D is the input channel size and (H, W) are the height and width of the input. The input values for D, H, W are 3, 512 and 512, respectively. The output of this module is in $\mathbb{R}^{B \times 512 \times (H/8) \times (W/8)}$.

2) *Fine-tuning Module*: The features extracted in Sec III-B.1 are fine-tuned for localization task in this module. Also, this module makes the features more suitable for the subsequent slot attention-based filtering. In the fine tuning module, the input is passed through a series of 2D convolutional layers, and is augmented with positional information channel-wise. The resultant features map is flattened along the spatial dimensions and fed into a linear layer for further processing. The positional encoding in this module is learnable with the encoding tensor of size $\mathbb{R}^{B \times (H/32) \times (W/32) \times 1024}$. The final output size of this module is $B \times (HW/1024) \times 1024$.

3) *Slot Attention-based Features Filtering*: Two types of distortions can considerably affect the localization performance. One results from the angular and range errors of the sensors, while the other from the foreground moving objects, e.g., bikes, buses, trucks and pedestrians in the outdoor scene. To minimize these distortions, state-of-the-art methods apply feature filtering techniques. For instance, Barnes *et al.* [22] designed a UNet-type architecture to predict distraction-free Radar odometry. Radarloc [8], AtLoc [25] and PointLoc [4] apply attention-based encoder and decoder techniques to filter out these noises. However, these methods fail to leverage object-centric features in the scene for this purpose, and they are also not easily tailored to multiple modalities. To this end, we design a unique slot attention [28] based feature filtering module to leverage object-centric features with Radar and image modalities. Slot attention is the key component of our module, whose architectural details are given in Fig. 1. Its primary function is to project the N input features to K output vectors, which represent slots - we chose $K = 20$ in our framework.

In this module, the slots are randomly initialized, and they improve iteratively during the training phase. The slots contest for input features through softmax-based attention throughout each iteration, and then use a gated recurrent unit (GRU) to update their representation. Similar to a transformer, “key”, “query”, and “value” vectors are used in the slot attention mechanism. Their computations is denoted by $k(\cdot)$, $q(\cdot)$, and $v(\cdot)$, respectively in the text to follow. These vectors are kept learnable in our slot attention to map the inputs and slots with a channel size D . The size of the affinity matrix for slot attention is $N \times K$ as compared to $N \times N$

for the multi-head attention (MHA) utilized in conventional transformer architectures, where $N \gg K$. This makes slot attention much more efficient than MHA. The slot attention output in $\mathbb{R}^{B \times 20 \times 1024}$ is passed through an MLP and a normal layer for feature refinement. The size of the feature vector generated here is $\in \mathbb{R}^{B \times 1024}$.

Concretely, a single iteration of the employed slot attention performs the following computation.

$$\alpha = \frac{1}{\sqrt{C}} k(input) \cdot q(slots)^\top \in \mathbb{R}^{N \times K}, \quad (1)$$

$$\Gamma_{i,j} = \frac{e^{\Lambda_{i,j}}}{\sum_s e^{\Lambda_{s,j}}}, \quad (2)$$

$$\beta = W^\top \cdot v(input) \in \mathbb{R}^{K \times D}, \quad (3)$$

$$W = \frac{\Gamma_{i,j}}{\sum_{s=1}^N \Gamma_{s,j}}. \quad (4)$$

In the above, α , Γ , and β respectively represent the attention coefficient matrix, normalized attention over the slots, and the slot update for further processing by the GRU. Finally, a GRU with as many hidden units as the dimensionality of the slots is used to update the slot. The update is based on a previous slot state and the signal β .

C. Modality Encoding

Inspired by the transformer’s positional encoding technique, we propose modality encoding and incorporate it in our network to identify the sensor modality. Unlike pre-fixed ‘cosine’ and ‘sine’ positional encoding, we learn the modality encoding for optimal performance. Since our framework supports three modalities, i.e., image, Radar and point cloud, we randomly initialize three modality encoding vectors with the same size as the feature vectors. These vectors are learned during the training phase and are added to the feature vectors to detect the sensor modalities. Experiments showed that this modality encoding works well during inference time.

D. Regression Module

The feature vectors encoded in the earlier stages of the model are passed to a Regression Module which is responsible to predict the 6DoF pose for localization. It comprises two common fully connected (FC) layers at the top, preceded by two parts of four FC layers - see Fig. 1. The network has a dicephalous architecture to precisely predict the translation and rotation parameters for the 6DoF pose. The channel output sizes in each division of the FC layers are 1024, 512, 256, 3. The initialization of the layers is set with Xavier_uniform distribution and ReLU activation is used. To reduce the variations between rotation and translation values, the rotation branch of this module is normalized.

E. Training Loss Computation

We utilize ℓ_1 -loss for rotation and translation, which we found more suitable for our network due to the various types of input data. We combine the rotation and translation losses to compute the loss for each input data with learnable balancing factors α and β , as shown in Eq. (5). During the training phase, we forward pass inputs from all six sensors

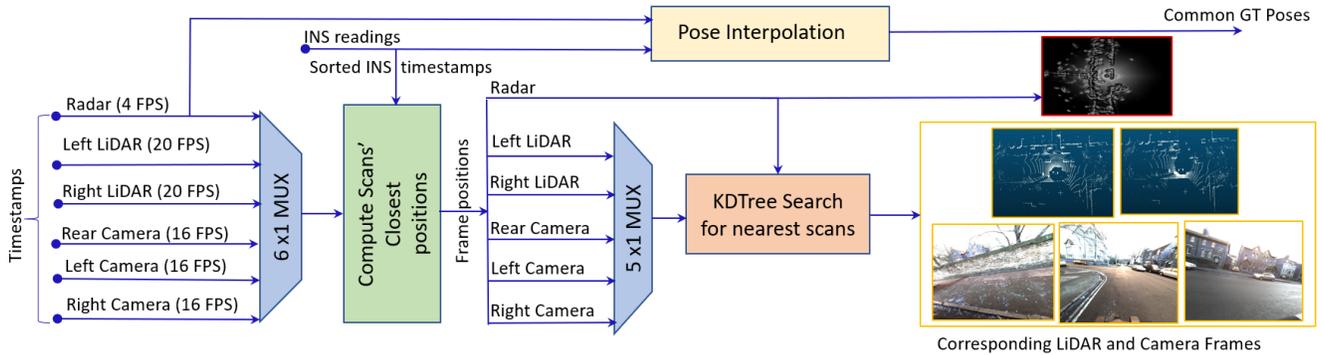


Fig. 4. Schematics for common ground-truth generation for the different types of input sensors namely, Radar, Camera and LiDAR which operate at 4, 16 and 20 frames per second respectively.

and combine the loss for each input to compute the net loss for the network, see Eq. (6). Finally, the net loss is back propagated in the network for optimization.

$$\mathcal{L} = \|t - t'\|e^\alpha + \alpha + \|r - r'\|e^\beta + \beta. \quad (5)$$

$$\mathcal{L}_{net} = \mathcal{L}_{L1} + \mathcal{L}_{L2} + \mathcal{L}_{C1} + \mathcal{L}_{C2} + \mathcal{L}_{C3} + \mathcal{L}_R. \quad (6)$$

In Eq. (5), t and t' indicate ground-truth and predicted translation, whereas r and r' denote the respective rotations. In Eq. (6), \mathcal{L}_{net} , \mathcal{L}_{L1} , \mathcal{L}_{L2} , \mathcal{L}_{C1} , \mathcal{L}_{C2} , \mathcal{L}_{C3} , \mathcal{L}_R are the net loss, LiDAR left, LiDAR right, camera left, camera right, camera rear, and Radar sensor losses, respectively.

IV. EXPERIMENTS

We evaluate the proposed method for the localization task on benchmark datasets and compare it with state-of-the-art techniques. To ensure a fair comparison, we choose popular existing methods based on the availability of source code by the original authors. We present results on three major localization datasets: Oxford Radar RobotCar [9], Apollo-SouthBay [10] and Perth-WA [3]. Our results establish the effectiveness of the proposed model for point cloud, image and Radar modality, both individually and collectively. Prior to presenting the results for each dataset, we discuss implementation details in the section below.

A. Implementation Details

To ensure fair benchmarking, we adopt uniform configurations for our method across the Oxford Radar RobotCar, Apollo-SouthBay and Perth-WA datasets. We apply batch sizes of 6 and 1 for training and testing, respectively. We use Adam optimizer with a learning rate of 0.0001 and weight decay 0.0005. At the outset, the model is trained on Oxford Radar RobotCar for multi-modalities and then fine-tuned on the Apollo-SouthBay dataset and Perth-WA dataset for point cloud modality. The model is trained for 40 epochs on all three datasets. For all experiments, NVIDIA GeForce RTX 3090 GPU with 24 GB memory is used. The experiments are conducted using PyTorch 1.8.0 on Ubuntu 18.04 OS.

B. Results on Oxford Radar RobotCar

Dataset details: The Oxford Radar RobotCar [9] is an extension of their previous dataset [32] that includes three different modalities: RGB camera images, Radar data, and point clouds from six sensors: left, right, and rear cameras; a

Navtech CTS350-XFMCW Radar scanner; and left and right Velodyne HDL-32E LiDAR. NovAtel SPAN-CPT ALIGN inertial (INS) and GPS navigation systems are used to collect the ground-truth poses for this dataset. It covers a total of 280km of urban area, including more than 30 sequences, each captured over 9km. This dataset is large and challenging for localization due to the presence of a variety of foreground objects, such as people and cars. We use the same training and test sequences as Radarloc [8] for our experiments on this dataset.

Common Ground-truth Generation for multi-sensors:

Our approach has the unique ability to leverage all three data modalities provided by the dataset. However, the frame rates for Radar (4Hz), LiDAR (20Hz), and camera (16Hz) data have a large disparity between them, which leads to timestamp misalignment between the modality sensors and the ground-truth. To generate a unified ground-truth for all the data sensors at a given time, we synchronize the Radar timestamps with the ground-truth poses using interpolation between GPS/INS measurements and the Radar timestamps. This step of ours provides ground-truth pose for each Radar frame. We then compute the position information for each frame for all modality sensors based on their timestamps. To acquire the corresponding frame for each remaining sensor, we search for the closest frame position corresponding to each Radar frame using the minimum Euclidean distance with KDTree search [33]. Each Radar frame and its corresponding closest searched frame share the same ground-truth. Missing GPS/INS data is handled by interpolating values from visual odometry data provided by the Radar RobotCar. In this way, we calculate single ground-truth pose for each frame of all modality sensors. The process is also illustrated in Fig. 4.

Performance: We present the experimental results in Table I where we use individual data modalities to compare with the approaches of the respective modalities. Due to the universal nature of our method, we are able to compare with Radar, camera, and LiDAR-based deep localization methods. Our technique outperforms all the existing methods by a considerable margin, which is clear from the respective sections of the table. For the point cloud modality, our approach surpasses the best performer PointLoc [4] by reducing the errors for

TABLE I

MEAN TRANSLATION (METERS) AND ROTATION (DEG.) ERRORS ON THE RADAR ROBOTCAR DATASET[9] COMPARED TO RADAR SLAM[29], ADAPTED ATLOC[25], RADARLOC[8], ATLOC[25], MAPNET[17], POSENET[5], DCP[30], VLAD[31] AND POINTLOC[4].

| Sq | Radar | | | | Camera | | | | LiDAR | | | |
|----|-------------------|---------------|-----------|--------------------|-------------------|------------|------------|--------------------|------------|---------------|------------|--------------------|
| | Radar SLAM | Adapted AtLoc | RadarLoc | UnLoc Ours | AtLoc | MapNet | PoseNet | UnLoc Ours | DCP | PointNet VLAD | PointLoc | UnLoc Ours |
| 6 | 49.8, 5.2° | 15.9, 4.2° | 8.4, 3.4° | 2.92, 1.26° | 15.4, 3.4° | 32.2, 5.4° | 51.1, 6.4° | 2.91, 1.27° | 18.5, 2.1° | 28.9, 5.2° | 14.4, 2.8° | 1.71, 0.68° |
| 7 | 24.7, 3.4° | 13.2, 3.8° | 5.1, 2.9° | 2.82, 1.68° | 39.7, 8.3° | 47.8, 5.4° | 80.3, 6.5° | 2.83, 1.69° | 02.8, 1.7° | 17.6, 3.9° | 08.5, 1.8° | 1.67, 0.68° |
| 8 | 26.1, 1.6° | 14.2, 2.9° | 6.6, 3.1° | 2.85, 1.21° | 31.7, 4.3° | 51.9, 7.7° | 111, 12.8° | 2.85, 1.20° | 16.4, 2.3° | 23.6, 5.9° | 09.5, 2.1° | 1.57, 0.61° |
| 9 | 39.8, 5.7° | 15.7, 3.2° | 6.5, 2.9° | 2.84, 1.49° | 47.1, 9.4° | 14.9, 2.8° | 45.5, 4.0° | 2.83, 1.50° | 13.6, 1.9° | 13.7, 2.6° | 11.5, 2.0° | 1.51, 0.35° |
| 10 | 17.8, 1.7° | 13.2, 1.9° | 5.3, 1.8° | 2.67, 2.04° | 10.4, 1.3° | - | - | 2.68, 2.05° | - | - | 08.4, 1.4° | 1.57, 0.51° |
| Av | 31.7, 3.5° | 14.4, 3.2° | 6.4, 2.8° | 2.82, 1.53° | 28.8, 5.3° | 36.7, 5.4° | 72.0, 7.4° | 2.82, 1.54° | 15.8, 2.1° | 20.8, 4.4° | 10.5, 2.0° | 1.61, 0.57° |

TABLE II

ABLATION STUDY ON TEST SEQUENCES OF OXFORD RADAR ROBOTCAR. MEAN TRANSLATION (METERS) AND ROTATION (DEG.) ERRORS ARE REPORTED FOR THE 3 MODALITIES SEPARATELY AND IN COMBINATIONS. L1/L2, C1/C2/C3 & R STAND FOR LiDAR LEFT/RIGHT, CAMERA LEFT/RIGHT/REAR & RADAR. WE REPORT RESULTS FOR C1 ONLY SINCE ALL THREE CAMERAS GIVE APPROXIMATELY SIMILAR ACCURACY.

| Seq | L1 | L2 | C1 | R | L1, C1, R | L2, C2, R | L1, L2, R | All |
|-----|-------------|-------------|-------------|-------------|-------------|-------------|-------------|--------------------|
| 6 | 1.75, 0.77° | 1.71, 0.68° | 2.91, 1.27° | 2.92, 1.26° | 1.68, 0.64° | 1.62, 0.59° | 1.50, 0.58° | 1.49, 0.57° |
| 7 | 1.69, 0.61° | 1.67, 0.68° | 2.83, 1.69° | 2.82, 1.68° | 1.63, 0.52° | 1.60, 0.61° | 1.47, 0.50° | 1.46, 0.49° |
| 8 | 1.58, 0.65° | 1.57, 0.61° | 2.85, 1.20° | 2.85, 1.21° | 1.52, 0.59° | 1.51, 0.57° | 1.38, 0.55° | 1.37, 0.54° |
| 9 | 1.55, 0.37° | 1.51, 0.35° | 2.83, 1.50° | 2.84, 1.49° | 1.49, 0.36° | 1.47, 0.35° | 1.34, 0.33° | 1.33, 0.32° |
| 10 | 1.52, 0.56° | 1.57, 0.51° | 2.68, 2.05° | 2.67, 2.04° | 1.47, 0.56° | 1.52, 0.48° | 1.38, 0.50° | 1.36, 0.48° |
| Avg | 1.62, 0.59° | 1.61, 0.57° | 2.82, 1.54° | 2.82, 1.53° | 1.56, 0.53° | 1.54, 0.52° | 1.42, 0.49° | 1.40, 0.48° |

translation and rotation nearly by $5\times$ and $3\times$, respectively. The results confirm that our encoding and sparse 3D convolutional modules are effective components for localization with point clouds. Table I also ascertains that our method is much more accurate than the camera and Radar-based methods. Our technique outperforms RadarLoc [8] by a more than $4\times$ error reduction in translation and rotation estimates. We conjecture that the strong performance of our approach has two main sources. Firstly, for each modality, our network is carefully designed with the state-of-the-art representation learning components. Secondly, our model is able to leverage the complementary information from different modalities during the learning stage to better train each individual modality network. For each modality, the inference stage is able to take guidance from learned positional encoding for optimal performance.

Ablation studies: To investigate the influence of different blocks of our technique, we conduct ablation studies and summarize the results in Table II. For the experiments, we kept the architecture for each block the same but turned on/off the image, Radar, and point cloud modality blocks to determine the impact on the localization results. First, we test our framework with a single modality and turned off the other two modality blocks. Table II reports the full results of these experiments in the first four columns, which can be compared with the results in Table I. Among the single modalities, our approach already achieves the best overall performance. To further analyse the performance of our network, we test different combinations of the blocks in our method using three sensors at a time. From Table II, it is clear that the localization performance of our technique improves by using different data modalities. In this case, the best performance is achieved with a combination of left Lidar, right Lidar and Radar, resulting in 1.42m and 0.49° errors for translation and

rotation, respectively. Finally, we utilized all six data streams in our technique. The last column of the table shows that this results in the overall best performance for our technique. This ascertains that each data modality is able to contribute to improve the performance, and that our network is able to leverage the complementary data information effectively.

C. Results on the Apollo-SouthBay Dataset

Dataset: ApolloSouthBay [10] is a comprehensive localization dataset collected in San Francisco, USA, utilizing an IMU-based system to record the ground-truth poses for the LiDAR frames. The dataset is captured in residential, urban, downtown area and highways. The dataset includes six routes: BaylandsToSeafood, ColumbiaPark, SanJoseDowntown, SunnyvaleBigLoop, Highway237 and MathildaAVE. All these routes provide separate training and test sets.

Performance: The outcomes of our experiments on this dataset are presented in Table III. For benchmarking, we fine-tune our model on the training sets and assess our model on all six routes of the test set. We employ RMSE as the evaluating metric by following [10] and compare our approach with the state-of-the-art localization methods, Levinson *et al.* [34], Wan *et al.* [35], L3-net [10] and Slice3D [3]. In the Table, we report average values on all six routes. Levinson *et al.*, Slice3D [3] and L3-net [10] are single modality methods, whereas Wan *et al.* [35] is a fusion model that integrates multiple sensors including a GPS. The results of compared methods are taken directly from the literature. Our method outperforms all techniques by achieving the lowest average errors across the six routes. We avoid reporting results of individual routes for brevity, however, note that our method achieves the best performance on each individual route as well.

TABLE III

RESULTS ON APOLLO-SOUTH BAY DATASET. THE VALUES REPRESENT RMSE FOR ROTATION (DEG) AND TRANSLATION (CM). OUR METHOD ACHIEVES THE LOWEST AVERAGE ERRORS FOR THREE ROUTES.

| Methods | <i>Yaw</i> | <i>Roll</i> | <i>Pitch</i> | <i>Rot</i> | <i>X</i> | <i>Y</i> | <i>Z</i> | <i>Trans</i> |
|-----------------------------|-------------|-------------|--------------|--------------|------------|------------|------------|--------------|
| Levinson <i>et al.</i> [34] | - | - | - | - | 11.9 | 9.3 | 4.6 | 8.9 |
| Wan <i>et al.</i> [35] | 3.8° | - | - | - | 5.0 | 3.6 | 2.6 | 3.7 |
| L3-net [36] | 1.6° | - | - | - | 5.0 | 3.6 | 2.7 | 3.8 |
| Slice3D [3] | 3.2° | 6.7° | 7.3° | 6.4° | 2.4 | 3.3 | 3.1 | 2.9 |
| UnLoc (Ours) | 2.8° | 1.3° | 1.4° | 1.83° | 3.1 | 2.7 | 2.6 | 2.8 |

TABLE IV

RESULTS ON PERTH-WA DATASET. THE VALUES REPRESENT THE ABSOLUTE MEAN ERROR FOR ROTATION (DEGREES) AND TRANSLATION (METERS). OUR METHOD HAS THE LEAST ERROR IN ALL CASES.

| Methods | <i>Yaw</i> | <i>Roll</i> | <i>Pitch</i> | <i>Rot</i> | <i>X</i> | <i>Y</i> | <i>Z</i> | <i>Trans</i> |
|---------------------------|--------------|--------------|--------------|--------------|-------------|--------------|-------------|--------------|
| PointLoc [4] | 0.26° | 1.96° | 0.15° | 0.75° | 29.70 | 37.49 | 7.80 | 25.00 |
| Slice3D (baseline) [3] | 0.32° | 2.42° | 0.27° | 1.00° | 14.20 | 17.05 | 8.50 | 13.25 |
| Slice3D (pre-trained) [3] | 0.17° | 1.52° | 0.10° | 0.59° | 6.26 | 06.55 | 2.86 | 5.23 |
| UnLoc (Ours) | 0.12° | 1.13° | 0.23° | 0.49° | 2.11 | 01.91 | 1.13 | 1.72 |

D. Results on Perth-WA dataset

Dataset: Perth-WA dataset [3] is captured in the Central Business District (CBD), Perth, Western Australia. The dataset comprises a LiDAR map of 4km² with 6DoF ground-truth per frame. The scenes include commercial structures, residential areas, food streets, complex routes, hospital buildings *etc.* The data is collected in three different two-hour sessions under various weather conditions. We apply the same split for training and testing sets as in [3]. The training set comprises 20K frames of sparse and dense point clouds, and another 2.2K frames are used as the test set. The dataset is available online on IEEE data portal [37]

Performance: We evaluate the performance of our approach against recent point cloud-based localization approaches: PointLoc [4], Slice3D baseline and pretrained models of [3], as shown in Table IV. To conduct this experiment, we fine-tune our framework on the training set and evaluate it on the test set. In line with Pointloc [4], we use the Mean Absolute Error of poses for analyzing the performance. Our localization approach outperforms all the methods for angular and translational mean error values. These results show that our proposed approach facilitates more effective point cloud feature learning, making it a preferred choice for outdoor localization using LiDAR frames.

V. CONCLUSION

This paper presents a novel localization framework for multi-sensors along with a deep neural network architecture that processes LiDAR, Radar and/or camera inputs on demand. The proposed network employs 3D sparse convolution and cylindrical partition to process LiDAR frames, and implements ResNet blocks with fine-tuning layers and a slot attention-based feature filtering module for the Radar and image modalities. It also introduces a novel learnable modality encoding technique to identify the type of input data modality. The network is trained on six inputs from three sensor types and can process either a single or multiple sensor inputs at inference. This makes our method robust to

sensor failure. Our method is useful for self-driving vehicles that need precise localization regardless of the weather conditions. We benchmark our method on three benchmark datasets and achieve state of the art results.

REFERENCES

- [1] A. Segal, D. Haehnel, and S. Thrun, "Generalized-icp." in *Robotics: science and systems*, vol. 2, no. 4. Seattle, WA, 2009, p. 435.
- [2] T. Stoyanov, M. Magnusson, H. Andreasson, and A. J. Lilienthal, "Fast and accurate scan registration through minimization of the distance between compact 3d ndt representations," *The International Journal of Robotics Research*, vol. 31, no. 12, pp. 1377–1393, 2012.
- [3] M. Ibrahim, N. Akhtar, S. Anwar, M. Wise, and A. Mian, "Slice transformer and self-supervised learning for 6dof localization in 3d point cloud maps," *arXiv preprint arXiv:2301.08957*, 2023.
- [4] W. Wang, B. Wang, P. Zhao, C. Chen, R. Clark, B. Yang, A. Markham, and N. Trigoni, "Pointloc: Deep pose regressor for lidar point cloud localization," *IEEE Sensors Journal*, vol. 22, no. 1, pp. 959–968, 2021.
- [5] A. Kendall, M. Grimes, and R. Cipolla, "Posenet: A convolutional network for real-time 6-dof camera relocalization," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 2938–2946.
- [6] A. Kendall and R. Cipolla, "Modelling uncertainty in deep learning for camera relocalization," in *2016 IEEE international conference on Robotics and Automation (ICRA)*. IEEE, 2016, pp. 4762–4769.
- [7] S. H. Cen and P. Newman, "Precise ego-motion estimation with millimeter-wave radar under diverse and challenging conditions," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2018, pp. 6045–6052.
- [8] W. Wang, P. P. de Gusmao, B. Yang, A. Markham, and N. Trigoni, "Radarloc: Learning to relocalize in fmcw radar," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2021, pp. 5809–5815.
- [9] D. Barnes, M. Gadd, P. Murcutt, P. Newman, and I. Posner, "The oxford radar robotcar dataset: A radar extension to the oxford robotcar dataset," in *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, Paris, 2020. [Online]. Available: <https://arxiv.org/abs/1909.01300>
- [10] W. Lu, Y. Zhou, G. Wan, S. Hou, and S. Song, "L3-net: Towards learning based lidar localization for autonomous driving," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6389–6398.
- [11] M. Elhousni and X. Huang, "A survey on 3d lidar localization for autonomous vehicles," in *2020 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2020, pp. 1879–1884.
- [12] X. Gao, Q. Wang, H. Gu, F. Zhang, G. Peng, Y. Si, and X. Li, "Fully automatic large-scale point cloud mapping for low-speed self-driving vehicles in unstructured environments," in *2021 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2021, pp. 881–888.

- [13] D. Kovalenko, M. Korobkin, and A. Minin, "Sensor aware lidar odometry," in *2019 European Conference on Mobile Robots (ECMR)*. IEEE, 2019, pp. 1–6.
- [14] J. Nubert, S. Khattak, and M. Hutter, "Self-supervised learning of lidar odometry for robotic applications," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2021, pp. 9601–9607.
- [15] R. Dubé, A. Cramariuc, D. Dugas, J. Nieto, R. Siegwart, and C. Cadena, "Segmap: 3d segment mapping using data-driven descriptors," *arXiv preprint arXiv:1804.09557*, 2018.
- [16] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," *arXiv preprint arXiv:1706.03762*, 2017.
- [17] S. Brahmabhatt, J. Gu, K. Kim, J. Hays, and J. Kautz, "Geometry-aware learning of maps for camera localization," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 2616–2625.
- [18] F. Walch, C. Hazirbas, L. Leal-Taixe, T. Sattler, S. Hilsenbeck, and D. Cremers, "Image-based localization using lstms for structured feature correlation," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 627–637.
- [19] M. Ding, Z. Wang, J. Sun, J. Shi, and P. Luo, "Camnet: Coarse-to-fine retrieval for camera re-localization," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 2871–2880.
- [20] V. Balntas, S. Li, and V. Prisacariu, "Relocnet: Continuous metric learning relocalisation using neural nets," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 751–767.
- [21] Z. Laskar, I. Melekhov, S. Kalia, and J. Kannala, "Camera relocalization by computing pairwise relative poses using convolutional neural network," in *Proceedings of the IEEE International Conference on Computer Vision Workshops*, 2017, pp. 929–938.
- [22] D. Barnes, R. Weston, and I. Posner, "Masking by moving: Learning distraction-free radar odometry from pose information," *arXiv preprint arXiv:1909.03752*, 2019.
- [23] D. Barnes and I. Posner, "Under the radar: Learning to predict robust keypoints for odometry estimation and metric localisation in radar," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2020, pp. 9484–9490.
- [24] X. Zhu, H. Zhou, T. Wang, F. Hong, Y. Ma, W. Li, H. Li, and D. Lin, "Cylindrical and asymmetrical 3d convolution networks for lidar segmentation," *arXiv preprint arXiv:2011.10033*, 2020.
- [25] B. Wang, C. Chen, C. X. Lu, P. Zhao, N. Trigoni, and A. Markham, "Atloc: Attention guided camera localization," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 06, 2020, pp. 10 393–10 401.
- [26] Z. Huang, Y. Xu, J. Shi, X. Zhou, H. Bao, and G. Zhang, "Prior guided dropout for robust visual localization in dynamic environments," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 2791–2800.
- [27] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [28] F. Locatello, D. Weissenborn, T. Unterthiner, A. Mahendran, G. Heigold, J. Uszkoreit, A. Dosovitskiy, and T. Kipf, "Object-centric learning with slot attention," *Advances in Neural Information Processing Systems*, vol. 33, pp. 11 525–11 538, 2020.
- [29] Z. Hong, Y. Petillot, and S. Wang, "Radarslam: Radar based large-scale slam in all weathers," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2020, pp. 5164–5170.
- [30] Y. Wang and J. M. Solomon, "Deep closest point: Learning representations for point cloud registration," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 3523–3532.
- [31] M. A. Uy and G. H. Lee, "Pointnetvlad: Deep point cloud based retrieval for large-scale place recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 4470–4479.
- [32] W. Maddern, G. Pascoe, C. Linegar, and P. Newman, "1 Year, 1000km: The Oxford RobotCar Dataset," *The International Journal of Robotics Research (IJRR)*, vol. 36, no. 1, pp. 3–15, 2017. [Online]. Available: <http://dx.doi.org/10.1177/0278364916679498>
- [33] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, "Scikit-learn: Machine learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.
- [34] J. Levinson and S. Thrun, "Robust vehicle localization in urban environments using probabilistic maps," in *2010 IEEE international conference on robotics and automation*. IEEE, 2010, pp. 4372–4378.
- [35] G. Wan, X. Yang, R. Cai, H. Li, Y. Zhou, H. Wang, and S. Song, "Robust and precise vehicle localization based on multi-sensor fusion in diverse city scenes," in *2018 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2018, pp. 4670–4677.
- [36] W. Lu, Y. Zhou, G. Wan, S. Hou, and S. Song, "L3-net: Towards learning based lidar localization for autonomous driving," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6389–6398.
- [37] M. Ibrahim, N. Akhtar, S. Anwar, M. Wise, and A. Mian, "Perth-wa localization dataset in 3d point cloud maps," 2023. [Online]. Available: <https://dx.doi.org/10.21227/s2p2-2e66>