

UFO: Uncertainty-aware LiDAR-image Fusion for Off-road Semantic Terrain Map Estimation

Ohn Kim*, Junwon Seo*, Seongyong Ahn, Chong Hui Kim

Abstract—Autonomous off-road navigation requires an accurate semantic understanding of the environment, often converted into a bird’s-eye view (BEV) representation for various downstream tasks. While learning-based methods have shown success in generating local semantic terrain maps directly from sensor data, their efficacy in off-road environments is hindered by challenges in accurately representing uncertain terrain features. This paper presents a learning-based fusion method for generating dense terrain classification maps in BEV. By performing LiDAR-image fusion at multiple scales, our approach enhances the accuracy of semantic maps generated from an RGB image and a single-sweep LiDAR scan. Utilizing uncertainty-aware pseudo-labels further enhances the network’s ability to learn reliably in off-road environments without requiring precise 3D annotations. By conducting thorough experiments using off-road driving datasets, we demonstrate that our method can improve accuracy in off-road terrains, validating its efficacy in facilitating reliable and safe autonomous navigation in challenging off-road settings.

I. INTRODUCTION

Autonomous navigation over complex and unstructured off-road terrains has become essential in developing a wide range of robotic applications, such as exploration, agriculture, and search and rescue. The effectiveness of off-road navigation is contingent upon the capability to accurately comprehend the relevant characteristics of the surrounding terrains related to navigational capability. Without prior knowledge of the environments, off-road navigation systems should be capable of examining terrain characteristics through onboard sensor measurements in real time [1]. The terrains can be classified semantically and translated into bird’s eye view representations, facilitating their integration into motion planning algorithms [2]. A semantic terrain classification map should be generated based on an accurate and comprehensive understanding of surrounding environments to ensure safe and effective navigation in off-road environments, which are characterized by rough and potentially hazardous terrains.

Producing a dependable semantic terrain classification map is challenging due to the distinctive characteristics of off-road environments, diverging significantly from indoor or structured settings, such as uncertain terrain boundaries and a wide range of terrain types [3], [4]. The highly variable terrain classes in off-road environments necessitate fine-grained labeling for comprehensive scene understanding [5]. Also,

This work was supported by the Agency for Defense Development Grant funded by the Korean Government in 2024.

The authors are with the Agency for Defense Development, Daejeon 34186, Republic of Korea {ohnkim.00, junwon.vision, seongyong.ahn, chonghui.chkim}@gmail.com

*These authors contributed equally to this work.

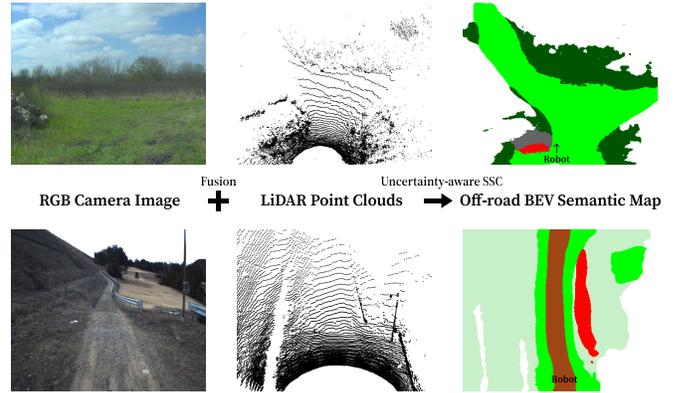


Fig. 1: Our method generates an off-road semantic terrain classification map in BEV from an RGB camera image and a single-sweep LiDAR point cloud in off-road environments. The classification accuracy is improved by combining complementary features extracted from the RGB image and the LiDAR point cloud. Moreover, our approach utilizes uncertainty-aware pseudo-labels to perform semantic scene completion (SSC), resulting in a dependable, dense semantic BEV map in diverse environments.

high intra-class variation of terrain appearances introduces potential unreliability in terrain classification outcomes [6]–[8]. Lastly, the complex geometry of off-road terrains makes mapping in a bird’s eye view with accurate geometry challenging due to the inability to estimate elevation accurately and adopt flat-ground assumptions [9].

These characteristics impose numerous restrictions on the reliable application of existing semantic terrain classification methods in various off-road environments. While it is expected to project pixel-wise segmentation into a Bird’s Eye View (BEV) map using LiDAR or stereo cameras [1], they produce sparse mapping by accumulating point-wise predictions from multiple time steps. Consequently, recent approaches leverage the concept of Semantic Scene Completion (SSC), employing learning-based inpainting to generate dense classification maps [10], [11]. However, these approaches encounter limitations due to the intrinsically sparse and less semantically rich features of LiDAR measurements. Some works explored deep learning-based models for terrain classification based solely on visual data, which utilize learning-based viewpoint transformations with 2D segmentation outcomes [12]–[16]. Nevertheless, without explicit range measurements, these methods produce geometrically inaccurate and unreliable outcomes, which might result in fatal failure during high-speed navigation [9]. Additionally, their dependence on specific data distributions of labeled datasets makes extending their application to diverse environments challenging.

This paper presents a terrain classification method that accurately estimates semantic maps in BEV in off-road scenarios, as shown in Fig. 1. Using a single LiDAR scan and an RGB image as input, the network generates dense and accurate terrain class maps in BEV by fusing information from multimodal sensor measurements. To ensure the applicability of the network in off-road scenes without relying on 3D annotations, pseudo-labels are generated through image-guided annotations. The reliability of the training with the pseudo-label is enhanced by uncertainty estimation, enabling the network to perform effectively in varied off-road environments. To validate the performance of our methodology, comprehensive experiments have been conducted using the publicly available RELIS-3D dataset [17]. Our method shows improved accuracy compared to single-modal methods by incorporating the multi-modal feature fusion methodology for semantic terrain map estimation.

II. RELATED WORKS

A. Robotics Mapping in Bird’s Eye View

Robotics mapping involves establishing a representation of a robot’s surroundings using noisy measurements as it moves through an environment [18]. The characteristics of these environments are assessed based on terrain features such as occupancy [19], traversability [20], or semantic class [21], [22]. These representations are commonly converted into the bird’s eye view due to their compatibility for integration with path planners in various robotic applications [23], [24]. To navigate efficiently and securely, the robot should be able to construct a map around itself online.

In the field of on-road navigation, camera-based methods widely employ viewpoint transformation learned by projecting pixel-wise features into BEV space [12]–[16], [25]. Nevertheless, the practicality of implementing these approaches in off-road settings is hampered by significant challenges, primarily due to real-time constraints and the absence of 3D terrain information. While other methods adopt range sensors such as LiDAR, they face challenges stemming from the sparsity of LiDAR returns despite its high precision of geometric information [2], [20]. Some methods perform semantic segmentation directly from raw sensor measurements, such as image or LiDAR, and then project the results onto BEV for the mapping [1]. However, the sparsity issues become more pronounced during high-speed navigation, where larger robot motion between LiDAR scans results in fewer depth measurements per unit area, compromising the reliability of the generated map.

Recent works propose learning-based approaches for predicting complete dense maps at a fixed size for off-road and unstructured environments to address these limitations [9], [26], [27]. Specifically, they leverage the concept of Semantic Scene Completion to generate a complete 3D scene from a single LiDAR scan [28], [29]. For instance, BEVNet achieves dense and accurate off-road terrain semantic classification based on SSC [11]. However, these methods often struggle to acquire accurate semantic predictions in off-road environments solely using LiDAR features. They also necessitate 3D

ground truth for consecutive scans, which poses a challenge in extending their applicability to off-road environments.

B. LiDAR-Image Fusion

While single-modal methods often face challenges in complex environments due to the inherent limitations of the input sensors, combining different modalities through sensor fusion has proven notable performance improvement in various applications [30]–[34]. LiDAR point clouds provide precise geometric information but only capture sparse data and lack texture information [35]. On the other hand, camera images can offer detailed and dense semantic information, while implicitly or explicitly inferred geometric information is prone to errors [9]. Hence, the combination of LiDAR and camera modalities is beneficial for performing terrain classification in off-road environments, as they can enhance each other’s capabilities.

Input-level fusion methods employed a BEV or spherical projection to project image logits or features into LiDAR space to improve the performance of LiDAR networks [36], [37]. The feature-level fusion methods aim to enhance feature representation by sharing information between the features of 2D and 3D backbones [30], [34], [38], [39]. Notably, the 2DPASS [40] effectively leverages rich semantic information from images by transferring knowledge across different modalities during the learning feature representations. It also acquires richer semantic and structural information through multi-scale feature fusion. To enhance the performance of models that generate semantic terrain classification maps for off-road scenarios, it is necessary to combine LiDAR and image modalities. This is because constructing precise maps for off-road scenes involves comprehension of both complicated geometry and rich semantics.

III. METHODS

This section details our proposed learning method for generating a dense off-road semantic terrain classification map in BEV. First, a method for creating a pseudo ground truth for creating a dense semantic map is proposed. Then, the network structure is presented, which can generate a dense semantic terrain classification map in the robot’s local frame from sensor measurements. LiDAR-image fusion is adopted to enhance prediction accuracy, and uncertainty-aware training is incorporated to increase the reliability of our method in various off-road settings.

A. Image-guided Pseudo Ground-truth Generation

While labeled datasets can create BEV ground truth in constrained environments, their acquisition cost and limited applicability to specific sensor configurations and class definitions pose challenges in off-road conditions. To ensure reliability in various off-road settings, we adopt a pseudo-label-based approach for generating BEV ground truth. This strategy alleviates constraints associated with the scarcity and expense of 3D labels, enabling training across a wide range of off-road scenarios. The overview of pseudo-labeling is depicted in Fig. 2

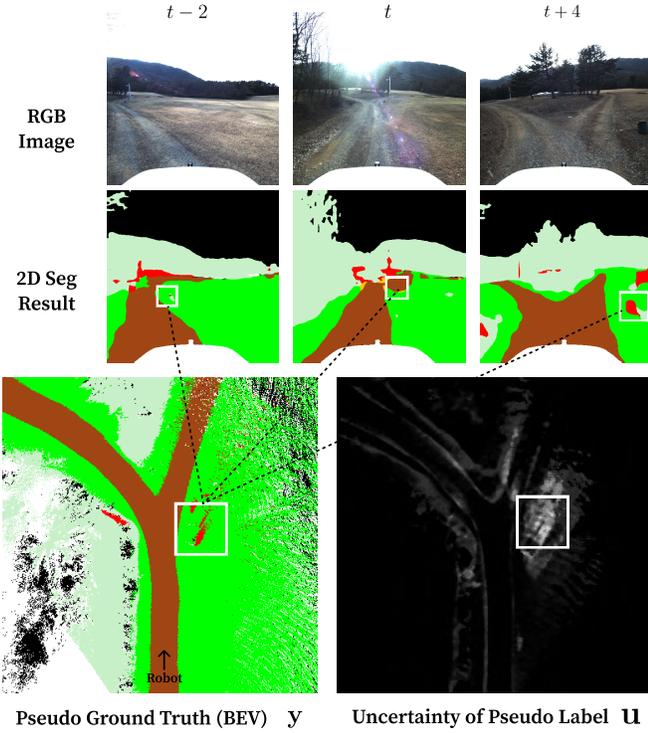


Fig. 2: Overview of image-guided pseudo-label generation. A pre-trained 2D image segmentation network derives semantic segmentation results from past and future images. These outcomes are aggregated using paired point clouds and then projected onto BEV grids to generate the pseudo-ground truth. Each grid determines a pseudo-label through the argmax operation, while its uncertainty is also quantified. Areas within the white box exhibit inconsistent semantic predictions across multiple timesteps, leading to higher uncertainties, depicted by brighter colors.

A pre-trained image segmentation network is utilized to produce the ground truth for semantic terrain classification. Image-based labels are beneficial in off-road scenes because they can accurately identify ambiguous boundaries and diverse classes, which are often characterized by rich details and textures. A dataset containing paired point clouds and RGB images can be easily obtained by navigating a robot equipped with LiDAR and a camera. The pseudo-labels of the BEV grid are then acquired by aggregating semantic segmentation predictions of the pre-trained model. Despite potential higher uncertainty in pseudo-labels, incorporating predictions from multiple sequential time steps during their generation along with uncertainty quantification can enhance their overall reliability.

For each paired LiDAR and RGB image, the inference results of 2D semantic segmentation are aggregated from the past 50 and future 100 data instances. Given the image segmentation results at timestamp t' , the prediction for each pixel is projected onto the BEV grid using the paired LiDAR point cloud. Given a i^{th} LiDAR point $\mathbf{P}^i \in \mathbb{R}^3$ at timestamp t' , the projection to pixel \mathbf{I}^i of each 3D point to a pixel in the image plane is determined based on camera parameters as follows:

$$\mathbf{I}^i = \mathbf{K} \cdot \mathbf{T} \cdot \mathbf{P}^i, \quad (1)$$

where $\mathbf{K} \in \mathbb{R}^{3 \times 4}$ and $\mathbf{T} \in \mathbb{R}^{4 \times 4}$ are the camera intrinsic and

extrinsic matrices respectively. The point is assigned a one-hot-encoded class prediction $\mathbf{S}^i \in \mathbb{R}^K$, where K is the total number of classes. The labeled points are then transformed into the reference robot frame at timestamp t as $\mathbf{P}_{t' \rightarrow t}^i$, based on the robot pose recovered by SLAM or odometry [41]. Multiple predictions are merged into a reference frame to produce dense BEV labels, and points from multiple timestamps are rasterized into BEV grid cells based on their x and y positions. A label for each grid G^j is determined from the points assigned to the grid, while some grids without assigned points are labeled as an unknown class. To avoid the aliasing of moving objects during aggregation, only points that belong to static object classes are aggregated, while points that belong to moving objects are aggregated only if they are from reference timestamp.

Class predictions of points assigned in the same grid are summed to calculate the grid class score, $\mathbf{c}^j \in \mathbb{R}^K$, where j is the index of the BEV grid. Then, the pseudo-label of a grid, \mathbf{y}^j , is determined through a majority vote of class predictions:

$$\mathbf{c}_k^j = \frac{1}{|G^j|} \sum_{i \text{ s.t. } \mathbf{P}_{t' \rightarrow t}^i \in G^j} \mathbf{S}_k^i, \quad (2)$$

$$\mathbf{y}^j = \arg \max_k \mathbf{c}_k^j. \quad (3)$$

Adopting multiple segmentation inference outcomes through majority voting introduces ensemble-like effects, effectively addressing potential inaccuracies in 2D segmentation results.

Additionally, the uncertainty of the pseudo-label of grid j , denoted as $\mathbf{u}^j \in [0, 1]$, is calculated to measure the reliability of the pseudo-label by measuring the consistency of segmentation for a grid over multiple timestamps:

$$\mathbf{u}^j = -\frac{1}{\log K} \sum_k \mathbf{c}_k^j \log(\mathbf{c}_k^j + \epsilon), \quad (4)$$

where $\epsilon = 1e^{-6}$ is used for stability. These uncertainty estimates can be leveraged during the network training to enhance the robustness of our semantic terrain classification map generation method [42].

B. BEV Semantic Fusion Network

Our network is trained to generate a dense top-view semantic classification map, utilizing a sparse frontal LiDAR point cloud and a paired RGB image. The pipeline of our method is depicted in Fig. 3.

LiDAR and image features are extracted in 3D voxel spaces using separate networks. The input LiDAR point cloud undergoes discretization into a (H, W, D) grid with a resolution of $0.1m \times 0.1m \times 0.2m$, and each point is further discretized into sparse voxels. In each voxel, a point is encoded as a 3-dimensional feature, comprising the offset from the voxel center $(\Delta x, \Delta y, \Delta z)$. Utilizing a simplified PointNet [43] architecture that includes a linear layer, Batch-Norm, and ReLU, each voxel of size $(N, 3)$ is transformed into sparse LiDAR voxel features of size C , where N is the maximum number of points per voxel. Simultaneously, image features are extracted from a pre-trained image backbone and

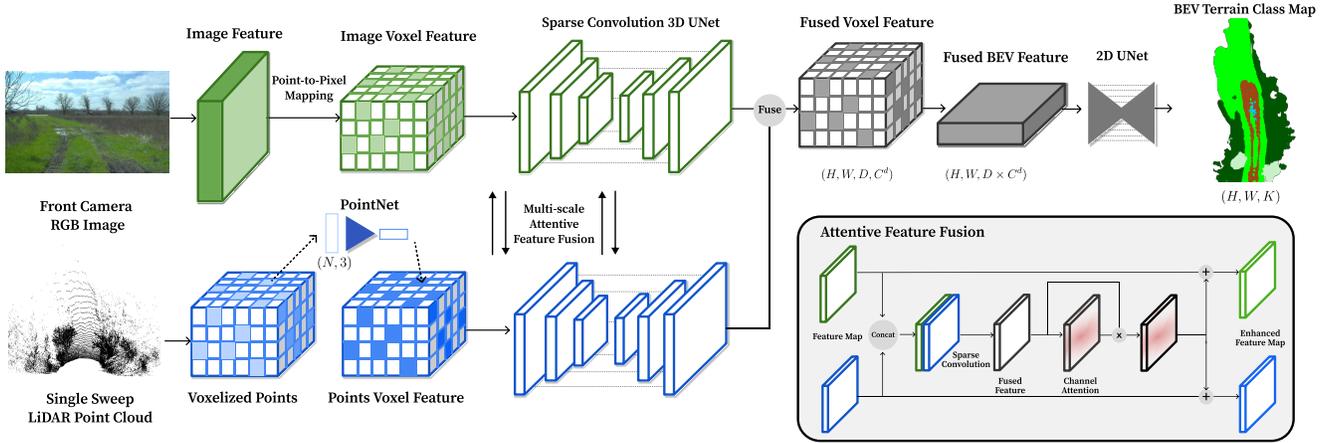


Fig. 3: High-level architecture of the proposed method. The network takes input from a single-sweep LiDAR point cloud and an RGB image captured by the front camera, producing a dense semantic terrain classification map in BEV. Extracted features from the image and point cloud, obtained through distinct encoders, are fused using Multi-scale Attentive Feature Fusion, integrated into the encoder of a 3D UNet of each modality. Subsequently, these fused features are passed to a 2D UNet to generate the dense semantic terrain classification map in BEV.

similarly converted into sparse image voxel features. The LiDAR points are projected onto the image plane via point-to-pixel mapping similar to Eq. 1, and the corresponding image features are propagated to the voxel to which the point belongs.

The input features from RGB and LiDAR are independently passed through separate 3D U-Net composed of sparse convolution layers to extract features for each modality. The 3D U-Net architecture employs a multi-level encoder-decoder structure, where each decoder layer is connected to the encoder of the same level through a residual skip connection. After each level of the sparse encoder, a multi-modal fusion block is integrated to promote feature fusion between LiDAR and image features. This facilitates the effective blending of rich semantic information from RGB with the geometric details derived from LiDAR points, ultimately leading to the generation of a precise semantic terrain classification map.

For feature fusion at each level, attentive fusion is employed to complement the features from each modality effectively, as shown in Fig. 3. Image and LiDAR features are concatenated channel-wise, followed by sparse convolution to generate a fused feature map. Channel attention is applied to the fused feature, emphasizing features that can strengthen the information from complementary sensors. The attended fused features are then added back to the original features of each modality, enabling a concentration on the more crucial information from each modality and distilling features from different modalities to enhance each feature further.

From the fused voxel features, a dense BEV terrain classification map is produced using a 2D convolution network. The voxel feature map, with dimensions (H, W, D, C^d) undergoes compression to yield a BEV feature map of dimensions $(H, W, C^d \times D)$ with empty grids initialized to zero. A U-Net-structured 2D convolution network is employed to generate a dense semantic map. This network progressively reduces the spatial size of features, capturing

higher-level semantic information, while the decoder part of the network upsamples feature maps to recover spatial information. Through this process, every grid feature is interpolated from sparse features to produce a dense terrain classification map of size (H, W, K) after a set of 1×1 convolutions in the segmentation head, which outputs logits for the classes.

C. Uncertainty-aware Terrain Classification

The model is optimized using the pseudo-label by employing uncertainty-weighted cross-entropy loss. To impose a lower weight on the ground truth with high uncertainty stemming from inconsistent pseudo-labels, the cross-entropy loss is calculated as follows:

$$\mathcal{L}_{\text{cls}} = - \sum_j \frac{\sum_k \mathbf{Y}_k^j \log \hat{\mathbf{Y}}_k^j}{1 + \mathbf{u}^j / \tau}, \quad (5)$$

where $\hat{\mathbf{Y}}^j \in \mathbb{R}^K$ is softmax probability for grid j , $\mathbf{Y}^j \in \mathbb{R}^K$ is one-hot vector representation of the label \mathbf{y}^j , and τ is the coefficient for controlling smoothness. By assigning weights inversely proportional to uncertainty in the cross-entropy term, the contribution of certain labels can be enhanced during optimization while that of uncertain labels is minimized. By utilizing this strategy, the network can be effectively trained without requiring manual annotations, while minimizing the negative impact on accuracy caused by pseudo-labeling.

IV. EXPERIMENTS

In this section, we validate the effectiveness of our method in enhancing terrain classification map generation performance. Through a quantitative and qualitative analysis, the results obtained from our method are compared with those of other existing approaches. This evaluation focuses on assessing the effectiveness of LiDAR-image fusion for terrain classification and confirming the validity of uncertainty-aware optimization for improving reliability.

TABLE I: Quantitative results on the RELIS-3D dataset [17]. Our method shows improved precision compared to those relying on a single modality.

Method	Acc [%]	mIoU [%]	 void	 dynamic	 static	 road	 grass	 dirt	 puddle	 rubble	 tree	 bush
<i>PyrOccNet</i> [13]	30.0	18.6	89.5	0.0	4.1	22.8	37.8	3.7	<u>4.7</u>	14.9	4.5	8.5
<i>TIM</i> [15]	30.5	17.6	88.8	0.0	1.2	9.6	36.6	6.1	3.2	13.5	7.2	9.4
<i>BEVNet</i> [11]	50.7	31.6	<u>91.6</u>	3.6	<u>11.1</u>	<u>40.1</u>	54.9	0.0	0.3	<u>35.0</u>	<u>54.3</u>	<u>25.0</u>
<i>Ours</i>	51.4	35.8	92.2	<u>0.9</u>	23.5	50.7	<u>54.3</u>	<u>4.6</u>	11.2	39.8	54.6	26.5

A. Dataset

We present our experimental results utilizing the publicly available off-road dataset, RELIS-3D [17]. This dataset contains RGB camera images and LiDAR point clouds, accompanied by point-wise semantic annotations and accurate robot poses recovered by SLAM [44]. We utilized sequences 0, 1, 2, and 4 for training and sequence 3 for evaluation. Note that this ensures the evaluation dataset contains distinct trajectory sequences not present in the training dataset. To address the class imbalance, some similar minor classes are grouped into a single category, resulting in a total of 10 classes, as shown in Tab. I.

B. Experimental Setup

1) *Implementation Details:* For all experiments, the input point cloud is cropped at $[(0, 51.2), (-25.6, 25.6), (-3, 5)]$ meters along the x, y, z axes, and a voxel grid size of $0.1m \times 0.1m \times 0.2m$ is used, resulting in dimensions $(H, W, D) = (512, 512, 40)$. The maximum number of points per voxel is set to $N = 10$, and the channels of voxel features are set to $C = 64$. For the 2D segmentation backbone, we adopt DeepLabV3 with ResNet50. The 3D U-Net sparse convolution network comprises encoder and decoder sections, each with four layers, and an attentive feature fusion block is attached to every encoder layer. The LiDAR stream encoder has one *SparseConv* and two *SubMConv* for each layer, while the image stream encoder has only one *SubMConv* for each layer. This design is because the pre-trained 2D image segmentation model has already extracted rich features for images. For the 2D U-Net, the encoder and decoder have four layers, respectively. Each downsampling and upsampling layer is connected with a skip connection, and each layer consists of multiple 3×3 convolutions, ReLU, and BatchNorm.

We train our model using the Adam optimizer, with a learning rate of $3e^{-4}$ and a batch size of 8 for 30 epochs. During training, point clouds are randomly augmented with vertical flipping, translations in the x and y axes within the range of $(-5, 5)$ meters, and rotations around the z -axis within the range of $(-\frac{\pi}{4}, \frac{\pi}{4})$ radians.

2) *Ground-truth Generation:* To establish the ground truth in the top view, we accumulate point clouds from consecutive LiDAR scans by transforming the points into the LiDAR coordinate frame of the current scan. The accumulation spans the past 50 scans and future 100 scans, with labels for each point determined by projecting the points into the semantic segmentation inference results of RGB images of the corresponding timestamp. Each point is then rasterized

into grids based on their x, y locations, and the semantic class and uncertainty for pseudo-labels are calculated for each grid. To mitigate the aliasing impacts of moving objects, only points belonging to static object classes are aggregated. For moving objects, only points from 3 sequential scans from the current scan are leveraged and given higher weight during the argmax operation.

3) *Comparison Methods:* Our method is compared to relevant approaches to evaluate our proposed method for generating the BEV semantic terrain classification map using LiDAR-image fusion. First, image-only methods that conduct view transformation to convert monocular images to BEV semantic maps are utilized. The Pyramid Occupancy Network (*PyrOccNet* [13]) employs a multiscale convolution network architecture to produce dense map representations directly from monocular images in BEV. Translating Images into Maps (*TIM* [15]) addresses generating BEV maps from images by solving sequences-to-sequences translation problems through an attention-based architecture. We also compare our method with a LiDAR-only approach that generates a dense semantic map in off-road, *BEVNet* [11]. Note that all models are evaluated solely with a LiDAR point cloud or RGB image of the current timestamp without using methods for temporal aggregation, such as recurrent neural networks, to objectively focus our evaluation on terrain map generation quality from sensor measurements.

4) *Evaluation Metrics:* Intersection over Union (IoU) is employed to assess the performance of semantic terrain classification in BEV. The IoU for class k is calculated as:

$$\text{IoU}_k = \frac{\text{TP}_k}{\text{TP}_k + \text{FP}_k + \text{FN}_k}, \quad (6)$$

where $\text{TP}_k, \text{FP}_k,$ and $\text{FN}_k,$ represent the number of true positive, false positive, and false negative predictions in grids, respectively. To evaluate the overall performance, the mean IoU (mIoU) is computed as:

$$\text{mIoU} = \frac{1}{K} \sum_{k=1}^K \text{IoU}_k, \quad (7)$$

where K is the total number of classes. Please note that unlabeled grids lacking ground truth labels due to the absence of accumulated points are excluded from the evaluation. For a comprehensive evaluation of our semantic terrain classification model, we additionally present overall prediction accuracy computed as:

$$\text{Accuracy} = \frac{\sum_{k=1}^K \text{TP}_k}{\sum_{k=1}^K (\text{TP}_k + \text{FP}_k)}. \quad (8)$$

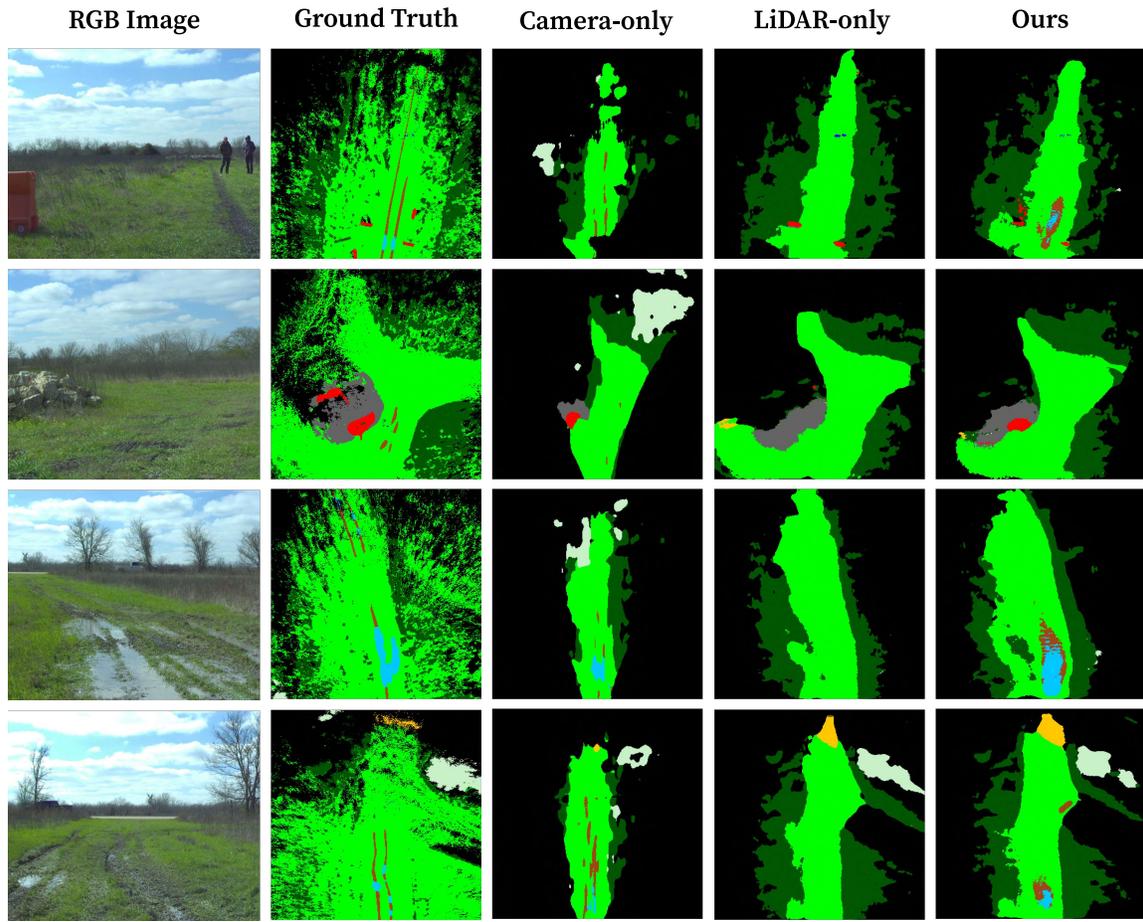


Fig. 4: Compared to other methods, *Ours* successfully predicted a semantic terrain map with the LiDAR-image fusion. The camera-only method excels at extracting semantic information from terrain but fails to represent a map with accurate geographical information in complex off-road environments. On the other hand, LiDAR-only methods accurately represent geographical information but are vulnerable to the semantic classification of terrain.

C. Experimental Results

The quantitative result for the RELLIS-3D dataset is presented in Tab. I. Our approach outperforms other methods specifically designed for structured conditions or that rely on a single modality. Image-only view-transformation methods perform poorly in off-road environments, especially for grass, road, tree, and bush classes containing intricate geometrical structures. LiDAR-based methods, on the other hand, exhibit more robust performance when predicting terrains with complex shapes. However, the accuracy of the LiDAR-only methods is inferior to ours, which improves reliability by utilizing uncertainty-aware optimization and sensor fusion. Specifically, LiDAR-only methods exhibit inferior performance in predicting classes requiring a higher-level understanding of textures, such as puddles and rubble.

Using LiDAR-image fusion, our method outperforms other baselines regarding mIoU and accuracy. Our method demonstrates improved IoU scores for classes challenging to distinguish solely with a single modality, such as dirt and puddles, validating the benefits of our fusion approach for accurately estimating the semantic properties of the surroundings. Qualitative results are presented in Fig. 4. Our method effectively produces semantic terrain maps with precise structures by

incorporating LiDAR features. Furthermore, our approach exhibits enhanced precision in estimating semantic terrain classes in off-road conditions. For example, our method accurately identifies puddles and dirt in BEV maps, which are challenging to capture solely with LiDAR point clouds.

D. Ablation Study

We present comprehensive ablation studies to examine the validity of each component of our methodology for estimating semantic terrain maps in off-road environments. Tab. II provides quantitative validation of the effectiveness of incorporating the LiDAR-image fusion component and uncertainty-aware optimization. The models without the fusion component are trained only using LiDAR features. Notably, incorporating multimodal fusion significantly improves performance, indicating strengthened features through fusion. While solely using uncertainty-weighted cross-entropy does not significantly enhance performance, it improves results when combined with the fusion methodology. This suggests that uncertainty-aware optimization effectively mitigates training uncertainty arising from imprecise labels and addresses the uncertainty of image features in pre-trained models.

TABLE II: Results of the ablation studies. Incorporating each module improves performance.

Module		RELLIS-3D	
Semantic Fusion	Uncertainty-aware Loss	mIoU	Acc
✗	✗	29.2	48.1
✗	✓	28.2	48.9
✓	✗	35.2	52.5
✓	✓	35.8	51.4

We then quantitatively evaluate the efficacy of our LiDAR-image fusion method, which fuses features of each modality on multiple scales with an attention mechanism. For comparisons, we train the network with modifications in fusion strategies. First, we conduct *Early fusion*, which simply concatenates image and LiDAR features before forwarding into the sparse convolution networks. Additionally, the network is trained without attentive fusion (*w.o. attention*), indicating that no channel attention is applied during the fusion step. Lastly, features are not fused in the multi-scale of the encoders but only in the last layer of the decoder (*w.o. multi-scale*), which is then forwarded to the 2D UNet for producing dense maps. To objectively assess the contribution of fusion methodologies, these models are compared with the model learned using our approach without uncertainty-aware loss (*Ours*).

TABLE III: Results of the ablation for LiDAR-image fusion.

Fusion Method	mIoU	Acc
<i>Early fusion</i>	29.7	47.7
<i>Feature fusion (w.o. attention)</i>	33.8	51.7
<i>Feature fusion (w.o. multi-scale)</i>	34.4	51.5
<i>Feature fusion (Ours)</i>	35.2	52.5

The experimental results for ablation studies for the fusion module are presented in Tab. III. It shows that our network design is effective in both datasets. While the early fusion methods, which simply concatenate the features, show improved results compared to the model that does not conduct fusion, it shows lower performance than feature-level fusion methods. This implies that more than simply decorating point features with image features is required to ensure the supplementation of the features effectively. Adopting attention during the feature fusion improves the performance, implying that the channel attention mechanism can facilitate overlapping two complementary features of each modality. Lastly, conducting feature fusions on multiple scales improves the results, suggesting the efficacy of fusing features at multiple resolutions, which aligns with other works that report the effectiveness of multi-scale fusions [16], [40].

V. CONCLUSION

This paper presents an approach to generating a dense terrain classification map in BEV. It can improve mapping accuracy through RGB-LiDAR fusion and enhance reliability

using uncertainty-aware pseudo-labeling. Utilizing a single LiDAR scan and an RGB image as input, the network employs attentive fusion at multiple scales to extract richer terrain features. Also, pseudo-labels are generated through image-guided annotations to enable the network to be learned without relying on 3D annotations. The training’s reliability with pseudo-labels is enhanced by uncertainty estimation, assigning lower weights to grids with high uncertainties. Evaluation using off-road driving datasets demonstrates the method’s efficacy for semantic terrain class map generation. The multimodal fusion approach proves advantageous in mitigating uncertainties associated with semantic understanding in challenging off-road scenes, where accurately assessing terrain properties poses greater difficulty than in structured environments.

Limitations and Future Works Although this method is highly effective in producing a precise semantic classification map of its surroundings in off-road environments, it is prone to generating overconfident predictions, which is a typical problem in map representations based on learning. This overconfidence arises from the explicit reliance on the network to generate the dense map. During interpolation from limited sensor measurements to infer semantic occupancy of grid cells, regions with occlusions may lead to overconfident predictions due to high uncertainties regarding terrain characteristics.

To enhance the applicability and reliability of our method in diverse real-world environments, we can leverage techniques for minimizing domain gaps, such as domain adaptation [45], [46]. Moreover, the potential of image data still needs to be fully exploited for BEV map generation. The method could benefit from incorporating view-transformation techniques, depth estimation, or leveraging recent successes in transformer architecture. Incorporating these features from dense image data could significantly enhance the generation of a dense map in BEV grids.

REFERENCES

- [1] D. Maturana, P.-W. Chou, M. Uenoyama, and S. Scherer, “Real-time semantic mapping for autonomous off-road navigation,” in *Field and Service Robotics: Results of the 11th International Conference*. Springer, 2018, pp. 335–350. [i](#), [ii](#)
- [2] J. Seo, T. Kim, K. Kwak, J. Min, and I. Shim, “Scate: A scalable framework for self-supervised traversability estimation in unstructured environments,” *IEEE Robotics and Automation Letters*, vol. 8, no. 2, pp. 888–895, 2023. [i](#), [ii](#)
- [3] X. Cai, M. Everett, J. Fink, and J. P. How, “Risk-aware off-road navigation via a learned speed distribution map,” in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2022, pp. 2931–2937. [i](#)
- [4] M. V. Gasparino, A. N. Sivakumar, Y. Liu, A. E. Velasquez, V. A. Higuti, J. Rogers, H. Tran, and G. Chowdhary, “Wayfast: Navigation with predictive traversability in the field,” *IEEE Robotics and Automation Letters*, vol. 7, no. 4, pp. 10 651–10 658, 2022. [i](#)
- [5] T. Guan, D. Kothandaraman, R. Chandra, A. J. Sathyamoorthy, K. Weerakoon, and D. Manocha, “Ga-nav: Efficient terrain segmentation for robot navigation in unstructured outdoor environments,” *IEEE Robotics and Automation Letters*, vol. 7, no. 3, pp. 8138–8145, 2022. [i](#)
- [6] M. Wigness, S. Eum, J. G. Rogers, D. Han, and H. Kwon, “A rugd dataset for autonomous navigation and visual perception in unstructured outdoor environments,” in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2019, pp. 5000–5007. [i](#)

- [7] J. Seo, S. Sim, and I. Shim, "Learning off-road terrain traversability with self-supervisions only," *IEEE Robotics and Automation Letters*, vol. 8, no. 8, pp. 4617–4624, 2023. [i](#)
- [8] X. Cai, M. Everett, L. Sharma, P. R. Osteen, and J. P. How, "Probabilistic traversability model for risk-aware motion planning in off-road environments," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2023, pp. 11 297–11 304. [i](#)
- [9] X. Meng, N. Hatch, A. Lambert, A. Li, N. Wagener, M. Schmittle, J. Lee, W. Yuan, Z. Chen, S. Deng *et al.*, "Terrainet: Visual modeling of complex terrain for high-speed, off-road navigation," *Robotics: Science and Systems (RSS)*, 2023. [i](#), [ii](#)
- [10] Y. Han, J. Banfi, and M. Campbell, "Planning paths through unknown space by imagining what lies therein," in *Conference on Robot Learning (CoRL)*. PMLR, 2021, pp. 905–914. [i](#)
- [11] A. Shaban, X. Meng, J. Lee, B. Boots, and D. Fox, "Semantic terrain classification for off-road autonomous driving," in *Conference on Robot Learning (CoRL)*, 2022, pp. 619–629. [i](#), [ii](#), [v](#)
- [12] J. Philion and S. Fidler, "Lift, splat, shoot: Encoding images from arbitrary camera rigs by implicitly unprojecting to 3d," in *European Conference on Computer Vision (ECCV)*, 2020. [i](#), [ii](#)
- [13] T. Roddick and R. Cipolla, "Predicting semantic map representations from images using pyramid occupancy networks," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 11 138–11 147. [i](#), [ii](#), [v](#)
- [14] Z. Li, W. Wang, H. Li, E. Xie, C. Sima, T. Lu, Y. Qiao, and J. Dai, "Bevformer: Learning bird's-eye-view representation from multi-camera images via spatiotemporal transformers," in *European Conference on Computer Vision (ECCV)*. Springer, 2022, pp. 1–18. [i](#), [ii](#)
- [15] A. Saha, O. Mendez, C. Russell, and R. Bowden, "Translating images into maps," in *International Conference on Robotics and Automation (ICRA)*. IEEE, 2022, pp. 9200–9206. [i](#), [ii](#), [v](#)
- [16] Z. Liu, H. Tang, A. Amini, X. Yang, H. Mao, D. L. Rus, and S. Han, "Bevfusion: Multi-task multi-sensor fusion with unified bird's-eye view representation," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2023, pp. 2774–2781. [i](#), [ii](#), [vii](#)
- [17] P. Jiang, P. Osteen, M. Wigness, and S. Saripalli, "Relis-3d dataset: Data, benchmarks and analysis," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2021, pp. 1110–1116. [ii](#), [v](#)
- [18] L. Gan, R. Zhang, J. W. Grizzle, R. M. Eustice, and M. Ghaffari, "Bayesian spatial kernel smoothing for scalable dense semantic mapping," *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 790–797, 2020. [ii](#)
- [19] K. Doherty, T. Shan, J. Wang, and B. Englot, "Learning-aided 3-d occupancy mapping with bayesian generalized kernel inference," *IEEE Transactions on Robotics*, vol. 35, no. 4, pp. 953–966, 2019. [ii](#)
- [20] J. Seo, T. Kim, S. Ahn, and K. Kwak, "Metaverse: Meta-learning traversability cost map for off-road navigation," *arXiv preprint arXiv:2307.13991*, 2023. [ii](#)
- [21] J. Wilson, Y. Fu, A. Zhang, J. Song, A. Capodiecchi, P. Jayakumar, K. Barton, and M. Ghaffari, "Convolutional bayesian kernel inference for 3d semantic mapping," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2023, pp. 8364–8370. [ii](#)
- [22] J. Wilson, J. Song, Y. Fu, A. Zhang, A. Capodiecchi, P. Jayakumar, K. Barton, and M. Ghaffari, "Motionsc: Data set and network for real-time semantic mapping in dynamic environments," *IEEE Robotics and Automation Letters*, vol. 7, no. 3, pp. 8439–8446, 2022. [ii](#)
- [23] M. Stölzle, T. Miki, L. Gerdes, M. Azkarate, and M. Hutter, "Reconstructing occluded elevation information in terrain maps with self-supervised learning," *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 1697–1704, 2022. [ii](#)
- [24] T. Miki, L. Wellhausen, R. Grandia, F. Jenelten, T. Homberger, and M. Hutter, "Elevation mapping for locomotion and navigation using gpu," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2022, pp. 2273–2280. [ii](#)
- [25] A. W. Harley, Z. Fang, J. Li, R. Ambrus, and K. Fragkiadaki, "Simplebev: What really matters for multi-sensor bev perception?" in *IEEE International Conference on Robotics and Automation (ICRA)*, 2023, pp. 2759–2765. [ii](#)
- [26] J. Fei, K. Peng, P. Heidenreich, F. Bieder, and C. Stiller, "Pillarsegnet: Pillar-based semantic grid map estimation using sparse lidar data," in *IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2021, pp. 838–844. [ii](#)
- [27] K. Peng, J. Fei, K. Yang, A. Roitberg, J. Zhang, F. Bieder, P. Heidenreich, C. Stiller, and R. Stiefelhagen, "Mass: Multi-attentional semantic segmentation of lidar data for dense top-view understanding," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 9, pp. 15 824–15 840, 2022. [ii](#)
- [28] R. Cheng, C. Agia, Y. Ren, X. Li, and L. Bingbing, "S3cnet: A sparse semantic scene completion network for lidar point clouds," in *Conference on Robot Learning (CoRL)*, 2021, pp. 2148–2161. [ii](#)
- [29] Z. Xia, Y. Liu, X. Li, X. Zhu, Y. Ma, Y. Li, Y. Hou, and Y. Qiao, "Scpnet: Semantic scene completion on point cloud," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023, pp. 17 642–17 651. [ii](#)
- [30] M. Liang, B. Yang, S. Wang, and R. Urtasun, "Deep continuous fusion for multi-sensor 3d object detection," in *European Conference on Computer Vision (ECCV)*, 2018, pp. 641–656. [ii](#)
- [31] Z. Chen, J. Zhang, and D. Tao, "Progressive lidar adaptation for road detection," *IEEE/CAA Journal of Automatica Sinica*, vol. 6, no. 3, pp. 693–702, 2019. [ii](#)
- [32] S. Pang, D. Morris, and H. Radha, "Clocs: Camera-lidar object candidates fusion for 3d object detection," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2020, pp. 10 386–10 393. [ii](#)
- [33] Z. Zhuang, R. Li, K. Jia, Q. Wang, Y. Li, and M. Tan, "Perception-aware multi-sensor fusion for 3d lidar semantic segmentation," in *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021, pp. 16 280–16 290. [ii](#)
- [34] Y. Li, A. W. Yu, T. Meng, B. Caine, J. Ngiam, D. Peng, J. Shen, Y. Lu, D. Zhou, Q. V. Le *et al.*, "Deepfusion: Lidar-camera deep fusion for multi-modal 3d object detection," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 17 182–17 191. [ii](#)
- [35] X. Bai, Z. Hu, X. Zhu, Q. Huang, Y. Chen, H. Fu, and C.-L. Tai, "Transfusion: Robust lidar-camera fusion for 3d object detection with transformers," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 1090–1099. [ii](#)
- [36] S. Vora, A. H. Lang, B. Helou, and O. Beijbom, "Pointpainting: Sequential fusion for 3d object detection," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 4604–4612. [ii](#)
- [37] C. Wang, C. Ma, M. Zhu, and X. Yang, "Pointaugmenting: Cross-modal augmentation for 3d object detection," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 11 794–11 803. [ii](#)
- [38] A. Piergiovanni, V. Casser, M. S. Ryoo, and A. Angelova, "4d-net for learned multi-modal alignment," in *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021, pp. 15 435–15 445. [ii](#)
- [39] D. Peng, Y. Lei, W. Li, P. Zhang, and Y. Guo, "Sparse-to-dense feature matching: Intra and inter domain cross-modal learning in domain adaptation for 3d semantic segmentation," in *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021, pp. 7108–7117. [ii](#)
- [40] X. Yan, J. Gao, C. Zheng, C. Zheng, R. Zhang, S. Cui, and Z. Li, "2dpass: 2d priors assisted semantic segmentation on lidar point clouds," in *European Conference on Computer Vision (ECCV)*, 2022, pp. 677–695. [ii](#), [vii](#)
- [41] T. Shan, B. Englot, D. Meyers, W. Wang, C. Ratti, and D. Rus, "Lio-sam: Tightly-coupled lidar inertial odometry via smoothing and mapping," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2020, pp. 5135–5142. [iii](#)
- [42] S. Ye, D. Chen, S. Han, and J. Liao, "Learning with noisy labels for robust point cloud segmentation," in *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021, pp. 6443–6452. [iii](#)
- [43] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, "Pointnet: Deep learning on point sets for 3d classification and segmentation," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 652–660. [iii](#)
- [44] W. Hess, D. Kohler, H. Rapp, and D. Andor, "Real-time loop closure in 2d lidar slam," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2016, pp. 1271–1278. [v](#)
- [45] S. Matsuzaki, H. Masuzawa, and J. Miura, "Multi-source soft pseudo-label learning with domain similarity-based weighting for semantic segmentation," *arXiv preprint arXiv:2303.00979*, 2023. [vii](#)
- [46] M. Jeon, J. Seo, and J. Min, "Da-raw: Domain adaptive object detection for real-world adverse weather conditions," *arXiv preprint arXiv:2309.08152*, 2023. [vii](#)