

TerrainNet: Visual Modeling of Complex Terrain for High-speed, Off-road Navigation

Xiangyun Meng¹, Nathan Hatch¹, Alexander Lambert¹, Anqi Li¹, Nolan Wagener², Matthew Schmitt¹, JoonHo Lee¹, Wentao Yuan¹, Zoey Chen¹, Samuel Deng¹, Greg Okopal¹, Dieter Fox¹, Byron Boots¹, Amirreza Shaban¹

¹University of Washington ²Georgia Institute of Technology

<https://sites.google.com/view/visual-terrain-modeling>

Abstract—Effective use of camera-based vision systems is essential for robust performance in autonomous off-road driving, particularly in the high-speed regime. Despite success in structured, on-road settings, current end-to-end approaches for scene prediction have yet to be successfully adapted for complex outdoor terrain. To this end, we present TerrainNet, a vision-based terrain perception system for *semantic* and *geometric* terrain prediction for aggressive, off-road navigation. The approach relies on several key insights and practical considerations for achieving reliable terrain modeling. The network includes a multi-headed output representation to capture fine- and coarse-grained terrain features necessary for estimating traversability. Accurate depth estimation is achieved using self-supervised depth completion with multi-view RGB and stereo inputs. Requirements for real-time performance and fast inference speeds are met using efficient, learned image feature projections. Furthermore, the model is trained on a large-scale, real-world off-road dataset collected across a variety of diverse outdoor environments. We show how TerrainNet can also be used for costmap prediction and provide a detailed framework for integration into a planning module. We demonstrate the performance of TerrainNet through extensive comparison to current state-of-the-art baselines for camera-only scene prediction. Finally, we showcase the effectiveness of integrating TerrainNet within a complete autonomous-driving stack by conducting a real-world vehicle test in a challenging off-road scenario.

I. INTRODUCTION

Autonomous robot navigation in off-road environments has seen a wide range of applications including search and rescue [50], agriculture [13], planetary exploration [48, 51], and defense [28]. Unlike indoor or on-road environments where traversable areas and non-traversable areas are clearly separated, off-road terrains exhibit a wide range of traversability that require a comprehensive understanding of the semantics and geometry of the terrain (Figure 1).

Current off-road navigation systems typically rely on LiDAR to obtain a 3D point cloud of the environment for semantic and geometric analysis [19, 33, 44, 45, 46]. While LiDAR sensors provide accurate spatial information, the resulting point cloud is rather sparse, making it tricky to build a *complete* map of the environment. Though point cloud aggregation can build such a map, it faces challenges when the vehicle travels at high speeds [20]. Finally, since LiDAR emits lasers into the environment, dust and snow can interfere with the measurement, and outside observers can detect the vehicle from the emitted lasers.



Fig. 1: High-speed driving in complex off-road environments requires joint reasoning of terrain semantics and geometry. Top row: a vehicle can drive at high-speed on a dirt road but has to be more cautious in snow due to wheel slipping. Bottom row: a vehicle needs to estimate terrain slopes and sizes of vegetation for safe planning and control.

Cameras, on the other hand, provide a number of benefits over LiDAR. Cameras provide high-resolution semantic and geometric information, stealth due to their passive sensing nature, are less affected by dust and snow, and are considerably cheaper. Hence, a camera-only off-road terrain perception system can potentially reduce the hardware cost, improve the system robustness at high speeds, and open up new possibilities for off-road navigation under extreme weather conditions and where stealth is desired.

Perhaps unsurprisingly, similar motivations have spurred recent major efforts of camera-only perception for *on-road* navigation [11, 22, 31, 35, 55]. This task mainly focuses on Bird’s Eye View (BEV) semantic segmentation to assess traffic conditions. One notable work is Lift-Splat-Shoot (LSS) [35]. The core of LSS consists of a “lift” operation that predicts a categorical distribution over depth for each pixel and a “splat” operation to fuse the image features and project them to the BEV space. LSS and related work are entirely data-driven, so they can predict complete maps and are more robust to sensor noise and projection errors. But their applicability to off-road perception is challenged by several barriers. First and foremost, they only predict a ground semantic BEV map without any 3D

terrain information that is critical for planning and control in off-road environments. Second, they are usually not optimized for real-time operation. For example, LSS predicts depth as a categorical distribution along the camera frustums to enable end-to-end learning, but this comes at a price: the size of the frustum features is large, creating a time and memory bottleneck, especially for field robots with limited hardware capabilities. Finally, to train such models, we need large-scale labeled terrain datasets. But, to the best of our knowledge, there are no such datasets for complex off-road terrains yet.

To this end, we design and implement TerrainNet, a real-time, camera-only terrain perception system that enables high-speed driving of a passenger-scale Polaris [2] vehicle on complex off-road terrains. We make several design choices and innovations to make TerrainNet suitable for off-road perception. First, TerrainNet supports multi-view RGB with *optional stereo depth* as inputs. Using stereo depth provides valuable geometric context that greatly improves prediction accuracy. Second, we enhance the predicted depth via an auxiliary loss on the output depth values. This extra supervision teaches the model to *correct* and *complete* the (potentially inaccurate) input stereo depth. This turns out to be critical for accurately estimating terrain geometry. To create ground-truth depth images, we build a complete map of the environment offline by aggregating LiDAR scans and removing outliers from the entire point cloud. Note that LiDAR is only used to create the training dataset and is not needed in deployment. Third, we make TerrainNet more than $5\times$ faster than LSS by lifting each image feature into a *single* 3D point and use a soft quantization technique in the “splat” step to keep the model end-to-end trainable. Lastly, TerrainNet predicts a *multi-layer* BEV map that captures both ground and overhanging terrain features.

TerrainNet is the first off-road, camera-only perception system for joint BEV semantic and geometric terrain mapping in a unified feed-forward model. To train and evaluate our model, we collect a new challenging large-scale, off-road dataset from different environments consisting of both *on-trail* and *off-trail* driving scenarios with *extreme elevation changes*. We believe our dataset better covers the diversity of the off-road driving challenges compared to RELLIS-3D [25], a publicly available dataset which is captured from a single environment and mainly consists of on-trail driving scenarios with limited elevation changes. We show that TerrainNet outperforms recent baselines in semantic and elevation estimation, while being much faster. Finally, we deploy TerrainNet inside a full navigation stack to have a Polaris vehicle traverse a 1.1 kilometer route over snow-covered hills.

II. RELATED WORK

On-road BEV perception. LiDAR and cameras are commonly used sensors in on-road autonomous driving perception systems [29, 35, 57], providing crucial information about surrounding objects and their semantics. Convolutional Neural Networks (CNNs) have shown exceptional performance in image [10] and point cloud [57] segmentation, and they have become the core of perception systems in on-road scene understanding. Although many of these systems rely on LiDAR

[29, 39], there has been increasing interest in using cameras due to their lower cost. In camera-based methods, a critical aspect is learning to project pixel-wise features to the BEV space. Lift-Splat-Shoot [35] adopts a *backward* projection scheme that lifts the image features using predicted depth and then splats the features into BEV space. SimpleBEV [21] and BEVFormer [31] perform *forward* projection from a set of grid points in the BEV space to retrieve their corresponding image features. Another line of work [11, 38, 40, 56] learns the projection with an attention mechanism. TerrainNet adopts a backward projection scheme since it makes full use of image pixels and does not assume a flat ground. Moreover, TerrainNet leverages stereo depth completion and soft-quantization for projection. This improves both the speed and accuracy of TerrainNet.

Off-road terrain modeling. Off-road terrains often exhibit large variations in ground elevations which can significantly affect terrain traversability. Hence, there has been a plethora of work on geometric terrain mapping. A frequently used representation is a 2.5D elevation map by aggregating point measurements from LiDARs or stereo cameras [16, 17, 18, 34, 46, 49]. In more complex environments where overhanging objects need to be considered, a voxel-based representation [5] or a multi-level surface map [49] are more effective in capturing detailed geometric information. In practice, obstacles or terrain discontinuities leave areas with missing values in the elevation map, leading to suboptimal motion planning. Inspired by data-driven image in-painting methods, recent works [41, 46] propose self-supervised learning to reconstruct the occluded area from an incomplete elevation map.

Besides geometric terrain modeling, semantics also play a critical role in traversability assessment [4, 26, 32, 33]. For instance, tall grass appears as obstacles yet is traversable, whereas large puddles are perceived as a flat surface but can be dangerous. Semantic segmentation is typically done in the image space [4, 42, 47] or the BEV space [33, 44]. Since planning is more convenient in the BEV space, it is a common practice to project the pixel-wise segmentation into a BEV cost map using LiDAR or stereo cameras [4, 33]. More recently, Shaban et al. [44] present an end-to-end trainable, recurrent CNN that builds a temporally consistent BEV semantic map from LiDAR. For a comprehensive literature review in traversability estimation for mobile robots, we refer readers to surveys [7, 43].

Existing off-road terrain perception approaches are usually designed for low-speed operations. Thus, they are able to gather dense sensor measurements to filter out the sensor noise to estimate a good terrain model. Some systems support high-speed operations but they are limited to on-road, flat terrains. In contrast, TerrainNet is designed for *high-speed* operations on *any terrain*. It achieves low latency by using a compact, multi-layer terrain representation [49], while at the same time maintaining accuracy at high speed with end-to-end BEV perception [35, 44]. Unlike existing systems that require complex algorithms and careful tuning to maintain real-time operations and consistent mapping, TerrainNet is a simple feed-forward neural network. It is also fast and can be improved as more training data is available.

Planning. To navigate a vehicle in off-road environments, a motion planner can leverage the perceived terrain features to plan safe and efficient trajectories. The terrain features are usually converted into costs for a planner to rank and assess the risk of trajectories [8, 9, 12, 15]. In our experiments, we illustrate how to use the MPPI planner [53] for motion planning with terrain features and robot capability considered.

III. OFF-ROAD TERRAIN MODELING

To enable a robot to drive safely and efficiently on off-road terrains, it is crucial to understand the traversability of its surroundings. Terrain traversability is the amount of cost or effort to traverse over a specific landscape. While many factors affect terrain traversability, we consider three primary factors: *semantics*, *geometry*, and *robot capability*.

Semantics. The semantics of terrain refers to the classes of objects (e.g., bush, rock, tree) or materials (e.g., dirt, sand, snow) occupying the terrain. Different semantic classes typically have different physical properties, such as friction and hardness, which can affect the capabilities of the vehicle. For example, since dirt can supply more friction than snow, a vehicle can drive faster on a dirt road than on snowy ground. Moreover, off-road vehicles have higher chassis and better suspension, so they can traverse over bushes and small rocks, albeit at lower speeds due to the increased resistance and bumpiness. Hence, the semantics of terrain encodes a rich spectrum of traversability.

Geometry. Off-road terrains are typically non-flat. A vehicle may not have enough power to climb a steep slope, and driving along a slide slope at high speed poses a significant risk of rolling over. Additionally, the geometry of objects also affects terrain traversability. For instance, a large bush is harder to traverse than a small bush. Hence, understanding the geometry of terrain is another important aspect of traversability assessment.

Robot capability. A vehicle’s physical and mechanical properties play another important role in terrain traversability. A bigger and more powerful vehicle can traverse over larger bushes or rocks than a smaller vehicle with less power. Since robot capability is an intrinsic property of the robot and is independent of terrain properties, we consider robot capability when designing the cost function later in Section VI-C. To some extent, robot capability is also considered when generating the training dataset, as described in Section V-B.

A. Multi-layer Terrain Representation

Since a ground vehicle traverses a 2D surface, it is convenient to use a gravity-aligned, 2D top-down grid map [14] to represent the terrain. Here we consider *local navigation*, where the map provides the vehicle with instantaneous information about its surroundings. Hence, we fix the size of the map and let the map “move” with the vehicle so that the vehicle stays at the center. This is commonly referred to as the *local map*. We do not consider building a global map in this work, though if required, we can leverage existing global SLAM algorithms to stitch the local maps together to obtain a global map.

The key question is what kind of terrain features to store in the top-down map. Previous work usually stores semantic classes [33, 44] or elevations [16, 46] for each grid cell. These representations have two key drawbacks. First, they do not model the hardness or porousness of the terrain, and hence they cannot capture the difference in traversability between a small and large bush. Second, they either do not consider overhanging objects or merge the semantic information of overhanging objects with ground objects. This would result in an inaccurate terrain model because the effect on traversability from overhanging objects is different from ground objects due to the geometry and the lack of gravity-induced force.

To address these issues, we extend the idea of MLS map [49] and propose a multi-layer terrain representation illustrated in Figure 2. It consists of two layers, a **ground** layer that captures terrain properties on the ground and a **ceiling** layer that models overhanging objects. For each layer, we model their semantic and geometric properties separately as follows:

Ground layer. For each map cell on the ground, we store the semantic probability distribution $\mathcal{C}_{\text{ground}}$ and elevation statistics $\mathcal{H}_{\text{ground}}$ of the terrain. The categorical distribution $\mathcal{C}_{\text{ground}} \in \mathbb{R}^K$ stores the relative proportions of the K semantic classes occupying each cell. We use the full distribution of semantic classes instead of a single class label to reduce the information loss caused by discretization (e.g., a map cell may contain both “dirt” and “bush” if it is at the boundary between dirt and a bush). The elevation statistics $\mathcal{H}_{\text{ground}}$ contains the minimum and maximum elevation values h_{\min} and h_{\max} on the ground. This allows the height of porous objects (such as grass and bushes) to be captured separately and resolves sharp elevation changes due to objects such as rocks and trees.

Ceiling layer. The ceiling layer models *overhanging* objects, such as canopies and tree branches. The semantic information $\mathcal{C}_{\text{ceiling}}$ is similar to $\mathcal{C}_{\text{ground}}$ but stores the semantic distribution of overhanging objects. The elevation information $\mathcal{H}_{\text{ceiling}}$ stores the height of the lowest overhanging point h_{ceiling} . If no overhanging points are present (e.g., on open terrains), we set $h_{\text{ceiling}} = h_{\min} + h_{\text{clearance}}$, where $h_{\text{clearance}}$ is a predefined constant of the desired vertical clearance.

This two-layer terrain representation captures a number of properties that are crucial for off-road navigation: the separate modeling of ground and overhanging semantics allows a robot to reason about semantic traversability more accurately, and the ground elevation statistics captures the hardness and sizes of ground objects (refer to Figure 2). These features can then be input to a cost function to build an accurate traversability cost map for planning (see Sec. VI-C).

B. Computing the Terrain Representation

This section illustrates a practical approach that leverages LiDAR sensors to build the proposed multi-layer terrain representation. This process is summarized in Figure 2. While our goal is real-time visual terrain modeling, we can benefit from accurate geometric sensors and extensive offline processing to produce high-quality ground-truth data for training that is otherwise challenging to obtain in real time.

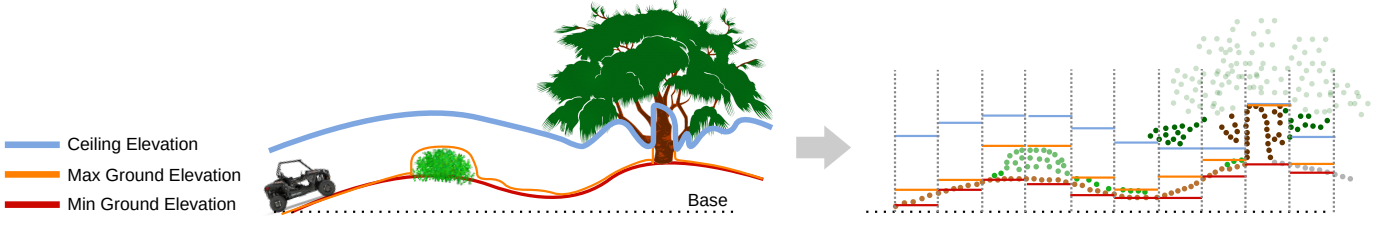


Fig. 2: **Left** Multi-layer terrain representation. The ground layer consists of the minimum and maximum ground elevations to capture the sizes of porous objects like bushes and also preserve sharp elevation changes after discretization. The ceiling layer represents the geometry of the overhanging objects (e.g., canopy) that might potentially block the path. **Right** In our practical implementation, we obtain a dense point cloud of the environment, discretize the x - y plane into 2D grid cells (shown in 1D for simpler visualization), and compute the elevation values by analyzing the point distribution in each cell. We ignore anything higher than the desired vertical clearance. For each layer and each cell, we count the number of semantic points for each class and compute a normalized histogram as the semantics of that cell.

Point cloud aggregation. Given recorded LiDAR scans, we use offline SLAM tools [23] to generate gravity-aligned poses for each scan and then aggregate them to obtain a globally aligned dense point cloud. We clean up the point cloud by removing points that are potentially outliers (see the supplementary for more details). Since this is done offline, we can thoroughly analyze the point cloud without worrying about the computational cost. Note that we could also use the stereo cameras for creating the dense point cloud, but we have found that stereo cameras have a shorter range and produce spurious depths on porous objects and the ground.

Divide and sort. We divide the x - y plane into 2D grid cells with a desired resolution and sort points that fall into a cell $C(i, j)$ by their z -coordinates in ascending order, giving us a sorted point set $\{(x_k, y_k, z_k) \in C(i, j)\}$ for each cell.

Compute elevations. We iterate over the points in each cell, starting from the lowest elevation. The minimum ground elevation of a cell is computed as

$$h_{\min} = \frac{1}{m} \sum_{k=1}^m z_k, \quad (1)$$

where m is a tuning parameter to smooth out the sensor noise. Cells without enough points are marked as unlabeled. Then, we compute the z -gaps of consecutive points as $\Delta z_i = z_{i+1} - z_i$. Letting z_{gap} denote a predefined constant corresponding to the minimum gap between the ground and the ceiling, if $\Delta z_i > z_{\text{gap}}$, then z_{i+1} is considered the lowest overhanging point, and we set:

$$h_{\max} = z_i \quad (2)$$

$$h_{\text{ceiling}} = z_{i+1}. \quad (3)$$

If no gap is found, then h_{\max} is set to the highest point in the cell, and $h_{\text{ceiling}} = h_{\min} + h_{\text{clearance}}$ where $h_{\text{clearance}}$ is the vehicle's vertical clearance.

Compute semantics. The ground and ceiling layers have separate semantic maps. Given a labeled point cloud (Section V-B), we assign each point to either the ground layer or the ceiling layer. Specifically, for a point (x, y, z) ,

$$(x, y, z) \in \begin{cases} \text{ground,} & \text{if } z \leq h_{\max}. \\ \text{ceiling,} & \text{if } h_{\max} < z < h_{\text{ceiling}}. \\ \text{ignored,} & \text{otherwise.} \end{cases} \quad (4)$$

After the partitioning, we compute a normalized histogram for each cell

$$[p_1, p_2, \dots, p_K] = \frac{1}{\sum_{k=1}^K n_k} [n_1, n_2, \dots, n_K], \quad (5)$$

where n_k is the number of points of class k in that cell.

IV. TERRAINNET

In the previous section, we introduced our terrain representation. When a vehicle is deployed in a new off-road environment, it has no knowledge of the terrain and thus must build the terrain representation from its onboard sensors in real time. In this section, we introduce our terrain inference engine *TerrainNet*, which predicts the terrain representation from its onboard RGB or RGB-D cameras.

A. Overview

TerrainNet is a single neural network that takes RGB or RGB-D images and jointly predicts terrain semantics and geometry. Figure 3 illustrates the overall pipeline. First, multiple RGB or RGB-D sensors are passed into the RGB-D backbone. The backbone produces two outputs for each camera: a dense and corrected depth image, and a 2D feature map. Each valid (i.e., neither self-hit nor sky) pixel in the depth image is back-projected to a 3D point using the corresponding camera intrinsics and extrinsics. We compute the terrain embedding for each 3D point that encodes semantic and geometric features. The per-point terrain embeddings are then down-projected and aggregated in the gravity-aligned BEV space with a fast, differentiable splat operation. The resulting BEV feature map is decoded into multiple semantic and elevation maps via an inpainting network. These maps are converted into a costmap by mapping the semantic categories and geometric features via a cost function that accounts for robot capability. The costmap can then be used by a local planner to find the best trajectory to reach a waypoint. In the following sections, we describe each component in detail.

B. Depth Completion and Correction

To build an accurate terrain representation, the most crucial aspect is to have a good spatial understanding of the environment. Modern stereo cameras can provide dense and accurate depth at close range. However, the error in depth estimation increases quadratically over distance, and stereo

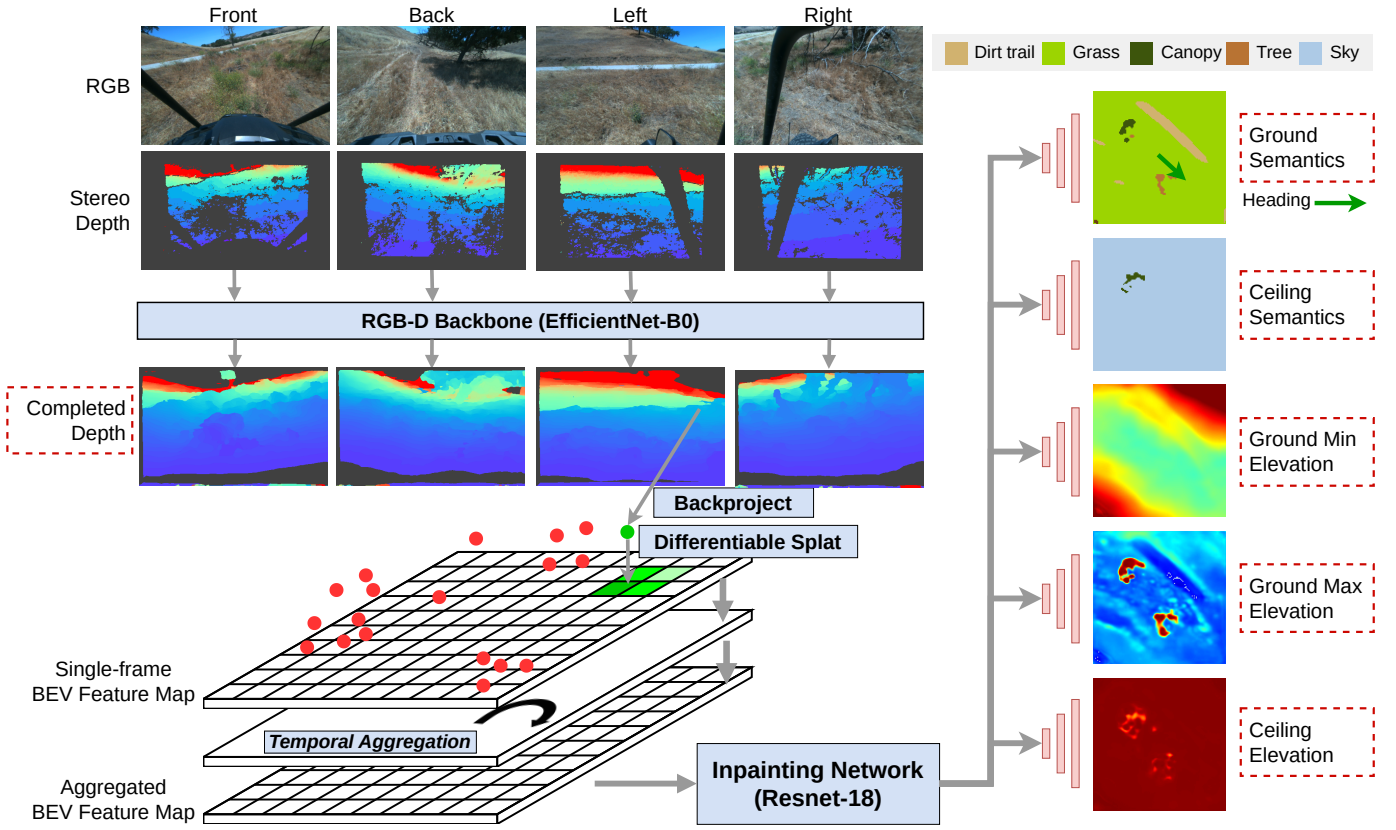


Fig. 3: **Architecture of TerrainNet.** RGB images with stereo depths are input to the RGB-D backbone to obtain completed depths with semantic features. Note that completed depths exclude pixels (colored in dark gray) that are potentially self-hits or sky. The per-pixel image features are back-projected to 3D to compute their terrain embeddings. The terrain embeddings are then splatted onto the BEV feature map with soft quantization. The feature map is temporally aggregated and finally passed to the multi-head inpainting network to predict a multi-layer terrain map. Red dotted boxes indicate where losses are applied during training.

matching works poorly for porous structures such as vegetation. Moreover, high-end stereo cameras [1] are equipped with high-sensitivity, low-distortion monochrome sensors for stereo matching, but they are not engineered for wide field-of-view semantic perception like RGB cameras. In Figure 3, we show the stereo depth output with the wide-angle RGB image on our hardware platform. The stereo depth only occupies a central area of the image, and there are many missing depth pixels due to failures in stereo matching. On the other hand, while there have been substantial advances in monocular depth estimation, they have only been shown to work well in structured environments where prior knowledge of object sizes and depth can be learned. Off-road terrains are often much less structured, and we find that solely using RGB information is inferior compared to a stereo system. To combine the rich semantic information of wide-angle RGB cameras and the incomplete stereo depth, we perform *Depth Completion*, which fills in missing depth pixels and corrects the errors in stereo depth in a data-driven fashion.

Given an RGB image and a stereo depth image (note that these may come from two different cameras), we first transform the stereo depth points into the RGB image to obtain an aligned RGB-D image I . Then, we pass the RGB-D image into a U-Net similar to Jaritz et al. [24] to get a dense depth map. The depth map does not have to be full-resolution, and we find

$8\times$ downsampling provides a good trade-off between speed and accuracy. We adopt a classification approach for depth completion because it captures depth discontinuity better than regression [37].

Given the completed depth image and the feature map, we use the camera intrinsics \mathbf{K} and extrinsics $[\mathbf{R}, \mathbf{t}]$ to compute the 3D location of each pixel in the gravity-aligned BEV space. Specifically, given a 2D pixel (u, v) with depth d , we transform it to the vehicle-centered, gravity-aligned BEV frame:

$$[x, y, z]^\top = \mathbf{R}\mathbf{K}^{-1}[u, v, d]^\top + \mathbf{t}. \quad (6)$$

Terrain embedding. For each 3D point, we use the corresponding image embedding as the semantic embedding f_{sem} . For elevation, we apply a multi-layer perceptron (MLP) on z to obtain the elevation embedding: $f_{\text{elev}} = \text{MLP}(z)$. We concatenate the semantic and elevation embeddings together and apply another MLP to obtain the per-point terrain embedding: $f = \text{MLP}(\text{concat}(f_{\text{sem}}, f_{\text{elev}}))$.

C. Fast and Differentiable BEV Projection

Given the 3D point set $\{(x, y, z, f)\}$, the next step is to project the points onto the BEV map by rounding each point into the nearest integer grid coordinates ($\text{round}(x/r), \text{round}(y/r)$), where r is the side length of each map cell. The main disadvantage of hard quantization is that we

can no longer correct the grid coordinates via back propagation due to the loss of gradient w.r.t grid coordinates. Inspired by related work in differentiable geometric learning [36], we adopt a **soft quantization** approach, where we “splat” each point feature into the 4 neighboring map cells with weights computed by bilinear interpolation. This allows the depth completion module to learn to adjust its depth prediction by minimizing the loss in the projected map.

Local feature aggregation. Since multiple points may fall into the same map cell, we compute the weighted average of the embeddings in each cell, where the weights are given by the soft quantization to create a grid-based BEV feature map. Let $S(i, j) = \{(f_1, w_1), \dots, (f_N, w_N)\}$ denote the set of point embeddings, and their corresponding weights within cell (i, j) , the feature embedding for the cell is computed as

$$f_{\text{cell}}(i, j) = \frac{1}{W} \sum_{k=1}^N w_k f_k, \quad (7)$$

where $W = \sum_{k=1}^N w_k$.

Temporal feature aggregation. Due to occlusion and errors in estimated depth, the map built from a single frame may not be sufficiently stable and complete. Hence we introduce an optional *Temporal Aggregation* (TA) layer that aggregates BEV feature maps over time. The TA layer is a single ConvGRU similar to that of Shaban et al. [44] but with one major difference: we use an orientation-stable odometry frame for aggregation instead of an ego-centric frame to better preserve fine-grained details. This is similar to the classical sliding window approach [33] where we shift the map and integrate the current sensor measurements, but we perform the shifting and aggregation on the feature map with a recurrent network.

D. Multi-head Terrain Inpainting Module

The terrain inpainting module decodes the BEV feature map into complete semantic and elevation maps. For semantic maps, it predicts a $H \times W \times K$ tensor to represent the probability distribution of the semantic classes, where H, W are the height and width of the map. For elevation maps, it predicts a $H \times W \times 1$ tensor to represent the corresponding elevation values. We adopt a U-Net as the inpainting module with a shared encoder and multiple convolutional decoding heads. Each head predicts a specific type of terrain map. This works better than a single decoding head with multiple channels because the semantic and elevation maps have different output spaces.

V. IMPLEMENTATION DETAILS

A. Hardware Platform

We collect our training and test data on a modified Polaris RZR vehicle [2] shown in Figure 1, which is capable of driving on off-road terrains at speeds up to 20 m/s (45 mph). The vehicle is equipped with 4 MultiSense stereo cameras [1] and 3 Velodyne 32-beam LiDAR sensors. LiDAR is only used for generating ground-truth with the method discussed in Section III-B and is not used during testing.

B. Dataset

We collected our training and validation data on the actual vehicle in several off-road environments with diverse appearances (Figure 4). We collected 5 sequences of manual driving data, totaling 20k frames. For each sequence, we ran Google Cartographer [23] to obtain gravity-aligned robot poses for building the ground-truth terrain maps and depth maps.

Create elevation maps. We follow the algorithm in Section III-B to generate the minimum ground, maximum ground, and ceiling elevation maps. For each time step t , we aggregate 300 LiDAR scans (running at 10 Hz) with a time range of $(t - 150, t + 149)$. Note that the vehicle is always at the origin. We set $m = 3$, $z_{\text{gap}} = 1$ m, and $h_{\text{clearance}} = 3$ m.

Create semantic maps. We ask human labelers to annotate 6k LiDAR scans. These scans are carefully selected such that they contained the relevant semantic objects. These annotated LiDAR scans are aggregated and projected to the BEV ground map or the BEV ceiling map using the algorithm in Section III-B. We identify 7 semantic categories: *dirt*, *dirt-trail*, *grass*, *bush*, *canopy*, *tree*, and *rock*. It is straightforward to extend the ontology and finetune the model if additional data is available, as we show in Section VI-C.

Create dense depth maps. We aggregate 50 LiDAR scans for each time step and project the point cloud to each RGB camera. For each pixel, we keep the smallest depth value.

C. Pseudo Labeling

To leverage unlabeled data, we train a LiDAR segmentation network to predict *pseudo labels* [30] for the unlabeled points in the training set. Specifically, we train Cylinder3D [57] on the labeled LiDAR points and predict the labels for all the unlabeled LiDAR points. We use the pseudo-labeled LiDAR points to build our ground-truth semantic maps for training. We *do not* apply pseudo labeling on the validation set. In Section VI-B2 we evaluate the effect of pseudo labeling.

D. Training

We split the dataset into 15k frames for training and 5k frames for validation by cutting each sequence into two non-overlapping train and validation segments. We manually labeled 4500 training frames and 1500 validation frames, and we pseudo-labeled the remaining training frames. All frames have elevation labels. The images are rectified and resized to 512×320 . We crop and warp the images such that they all have the same intrinsic parameters. We train TerrainNet on 4 Nvidia A40 GPUs with a batch size of 8 and the Adam optimizer [27]. All models (including applicable baselines) are trained for 100k iterations. The map size is $51.2 \text{ m} \times 51.2 \text{ m}$ with a resolution of 0.4 m. There is no limit on the elevation range.

Loss. We first pre-train the RGB-D backbone on the depth completion task. We discretize the depth into 128 bins with a max range of 25.6 m and apply a cross-entropy loss. Then, we train the whole network end-to-end. For the semantic maps, we use an unweighted cross-entropy loss. For the elevation maps, we use a Smooth-L1 loss [3] with transition parameter $\beta = 0.2$. Instead of predicting h_{max} and h_{ceiling} directly, we

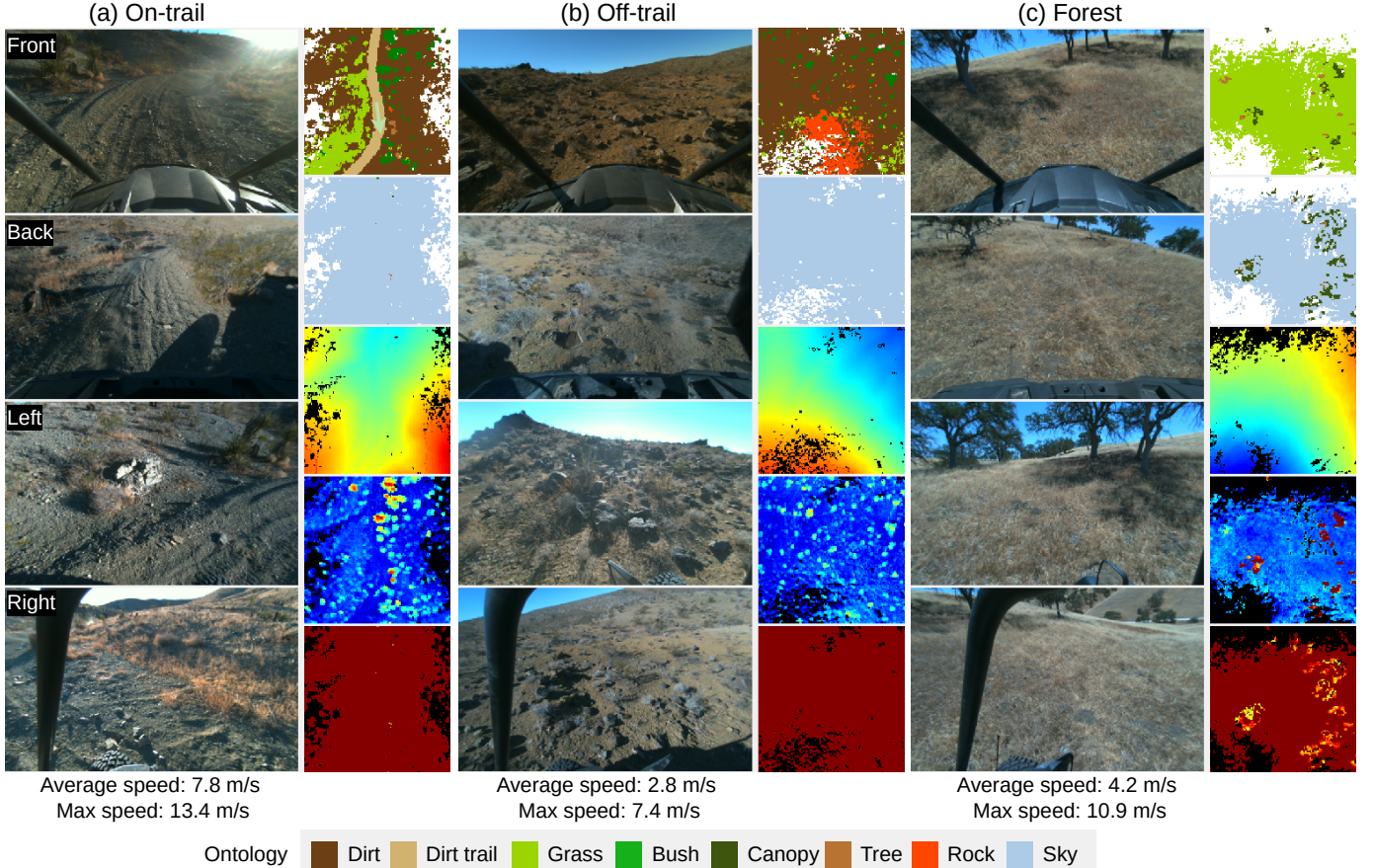


Fig. 4: Our dataset consists of diverse off-road environments. **a)** vehicle driving at high-speed on a trail with bushes, grass, rocks, and Joshua trees on the two sides; **b)** vehicle driving off-trail on a hilly terrain with scattered grass, bushes, rocks and Joshua trees; **c)** vehicle driving on a hilly terrain with tall trees and overhanging canopy. For each environment, we also show (from top to bottom) the ground semantics, ceiling semantics, min ground elevation, max ground elevation, and ceiling elevation. For the ceiling semantic map, we use the *sky* class to mark areas where no overhanging objects are present. The max ground and ceiling elevation maps are *offsets* to the ground elevation.

predict their *offsets* to h_{\min} . The total loss is $L = L_{\text{semantic}} + 0.1L_{\text{elevation}} + 0.1L_{\text{depth}}$.

VI. EXPERIMENTS

A. Comparison with Existing BEV Perception Methods

We quantitatively evaluate TerrainNet in terms of its accuracy and speed. We divide the experiments into a few sections, each focusing on a specific aspect.

1) *Baselines*: There is a vast amount of literature on on-road BEV perception. Many of them use heavyweight neural nets, which are unsuitable for running on a compute-limited vehicle. To this end, we focus on three main BEV perception paradigms and pick the most representative baseline for comparison:

Segmentation and Projection. We build a classical pipeline [19, 33] where we first perform image segmentation and then project the pixel labels using the stereo depths. Since we only label the LiDAR points, we project the labeled LiDAR points to each image as the ground-truth segmentation mask. We train a state-of-the-art image segmentation network SegFormer-B1 [54] on our dataset. Then, we use the same pipeline for generating the ground-truth maps to build a dense point cloud and project it down to the BEV map. We perform nearest neighbor inpainting to fill in the unknown space.

Lift-Splat-Shoot. LSS [35] implicitly learns to predict a per-pixel depth map and splats the features along each depth ray, weighted by the depth distribution. Since vanilla LSS cannot predict elevation, we modified LSS to compute the terrain embedding in the same way as TerrainNet. The main difference between LSS and TerrainNet is that TerrainNet has explicit depth supervision and does a “one-hot” splat of terrain embeddings onto the map.

SimpleBEV. SimpleBEV [21] performs forward projection with a uniform splat of image features without considering per-pixel depth. While it has been shown to outperform LSS and many other strong methods for on-road driving, it is not clear if the same holds for *off-road* driving where the terrain is not flat. We use a grid resolution of $128 \times 128 \times 32$ to cover a 3D volume of $51.2 \text{ m} \times 51.2 \text{ m} \times 51.2 \text{ m}$. (Using a larger z -resolution would make the memory consumption and training time prohibitive.)

We evaluate two versions of our system: the single-frame TerrainNet and TerrainNet-TA (i.e., with a temporal aggregation layer). For a fair comparison, LSS and SimpleBEV use the same image backbone and inpainting net as TerrainNet with the same hyperparameters for training. We also train a RGB-only version of TerrainNet, LSS, and SimpleBEV by removing the

	Method	Dirt	Dirt-trail	Grass	Bush	Canopy	Tree	Rock	mIoU \uparrow	Elevation (MAE) \downarrow			Time (ms) \downarrow
RGB	LSS	0.726	0.621	0.921	0.241	0.042	0.081	0.141	0.396	0.375	0.323	0.015	202
	SimpleBEV	0.727	0.614	0.922	0.245	0.040	0.041	0.106	0.385	0.417	0.322	0.018	110
	TerrainNet	0.730	0.617	0.919	0.257	0.071	0.081	0.177	0.407	0.361	0.315	0.018	25
RGB+Stereo	LSS	0.762	0.656	0.926	0.353	0.096	0.141	0.235	0.453	0.257	0.286	0.016	207
	LSS one-hot depth	0.751	0.642	0.919	0.308	0.094	0.098	0.202	0.431	1.028	0.379	0.019	28*
	SimpleBEV	0.754	0.647	0.923	0.320	0.076	0.070	0.140	0.419	0.286	0.299	0.015	113
	TerrainNet	0.765	0.666	0.926	0.380	0.125	0.145	0.274	0.469	0.244	0.277	0.015	28
	Seg & Proj	0.664	0.671	0.784	0.234	0.109	0.081	0.168	0.387	0.559	0.408	0.034	-
	TerrainNet-TA	0.796	0.679	0.928	0.513	0.161	0.192	0.276	0.506	0.243	0.240	0.024	29

TABLE I: For semantics, we report the per-class IoU and mean IoU. For elevation, we report the Mean Absolute Error (MAE) for min. ground, max. ground, and ceiling elevation, respectively. We report the inference time on an RTX 3090. The inference time for LSS one-hot depth is an estimation. For Seg&Proj we do not report its inference time since it depends on the implementation of temporal aggregation.

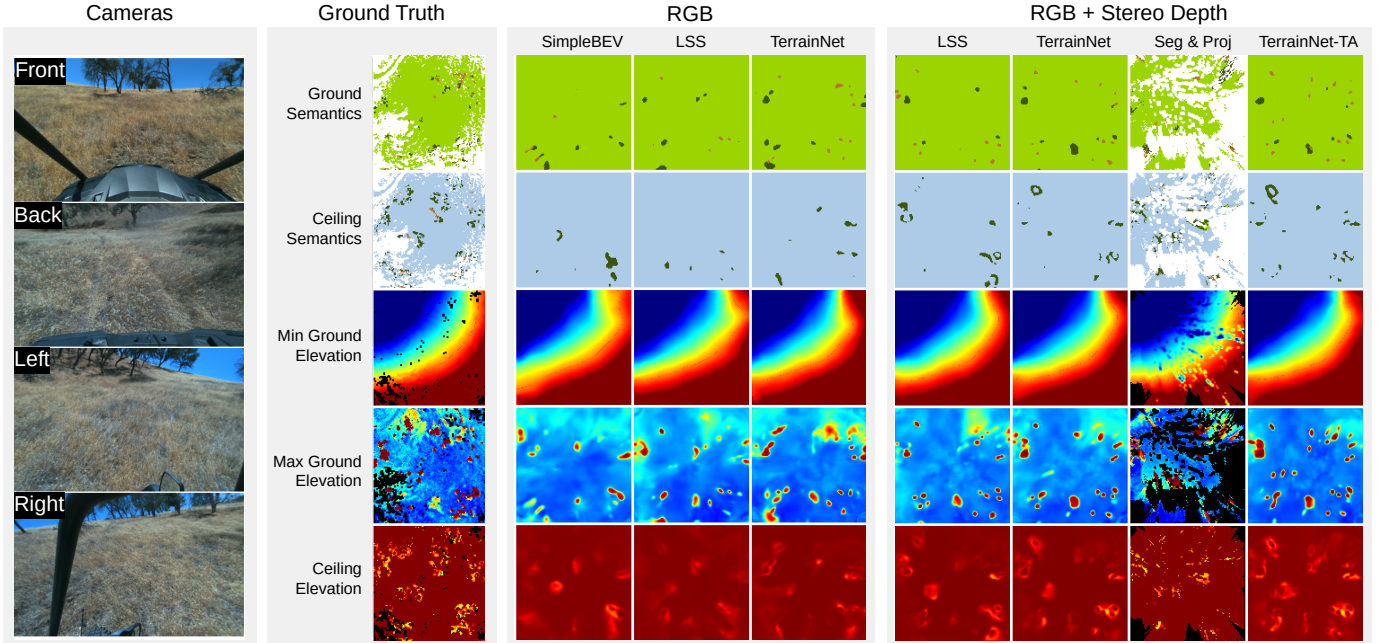


Fig. 5: Qualitative comparison between TerrainNet and the baselines in one of the forest environments (best viewed when zoomed in). In general, RGB-D models see more trees. Seg&Proj shows many artifacts due to errors in stereo matching on the ground. TerrainNet-TA shows the highest fidelity. See the supplementary material for more examples.

input depth.

2) *Terrain Modeling Accuracy*: In Table I, we compare TerrainNet with the baselines in terms of their accuracy in modeling the terrain. For semantics, we use the standard IoU metric. For elevation, we use the per-pixel mean absolute error.

TerrainNet surpasses all the baselines for RGB and RGB-D inputs. The fact that TerrainNet outperforms LSS shows that the explicit learning of per-pixel depth is beneficial. While it is possible to project a single depth for LSS during inference, the results are much worse, as shown in the *LSS one-hot depth* baseline. SimpleBEV and Seg&Proj perform worse than TerrainNet and LSS. Figure 5 presents a qualitative comparison and Figure 6 highlights a particular frame. The RGB-D models preserve more semantic details, with TerrainNet-TA being closest to the ground truth. The artifacts in the Seg&Proj baseline are due to errors in stereo matching. For more qualitative examples, please see the supplementary material.

In general, large and more frequent classes (*dirt*, *dirt-trail*, *grass* and *bush*) are easier to predict than small, less frequent objects (*tree* and *rock*). This is expected since small objects are harder to localize, especially when they are far. *Canopy* is also harder to localize due to its overhanging nature and

occlusion. In terms of elevation, TerrainNet-TA does not show notable improvement in the ground elevation than TerrainNet. We hypothesize that it is due to the recurrent layer not being able to track the vertical movement of the vehicle well. We leave it as future work to improve this aspect.

Stereo depth provides significant gains. Models with stereo depths as additional input outperform their RGB-only counterparts by a large margin. For off-road environments, objects are randomly scattered, and they can have similar appearances but different sizes (e.g., small vs. large bushes). Thus, it is hard to estimate the depth from color information. This suggests that for off-road terrain perception, it is preferred to equip the robot with stereo cameras instead of monocular cameras. The performance improvement is usually worth the extra cost.

Forward vs. backward projection. SimpleBEV performs on par with other methods on *dirt*, *dirt-trail*, and *grass*. These semantic classes cover a large area of the ground, so accurate depth estimation is not essential. But for objects such as *bush*, *canopy*, *tree*, and *rock*, SimpleBEV usually underperforms. SimpleBEV also performs poorly on elevation estimation, likely

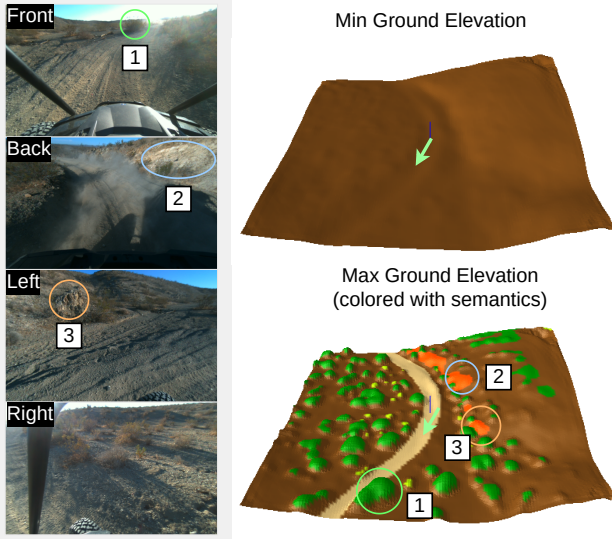


Fig. 6: 3D visualization of TerrainNet-TA prediction on a high-speed on-trail sequence. The green arrow indicates the vehicle’s heading. The Min Ground Elevation map shows a smooth ground support with porous objects such as vegetation removed. The Max Ground Elevation map contains the protruding vegetation on the ground. Note the bush far ahead (1) and rocks on the hills (2, 3) are well captured.

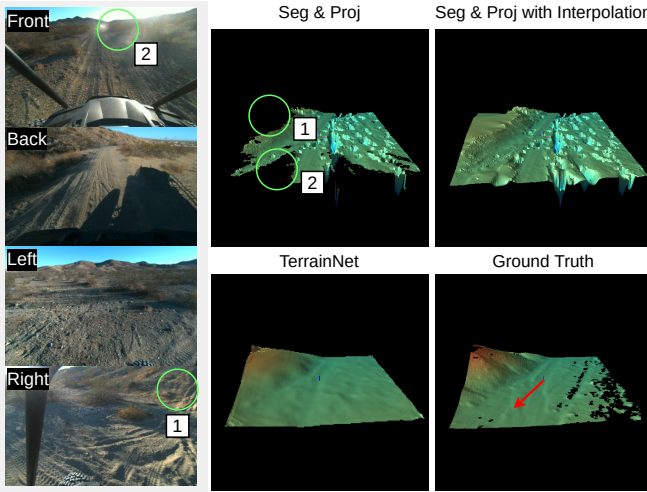


Fig. 7: Comparing TerrainNet to Seg&Proj with interpolation. TerrainNet is able to predict the non-visible hill crest (area 1) and the occluded flat ground (area 2) behind dense vegetation, whereas interpolating aggregated stereo depths predicts a flat hill and a bumpy ground. TerrainNet is also robust to artefacts in stereo depths.

due to the low z -resolution of the grid. Increasing the grid resolution would linearly increase computation and memory consumption, which is not always feasible.

Learning-based inpainting is better than interpolation. In Figure 7, we compare the predicted ground elevation map from Seg&Proj with Navier-Stokes interpolation [6] and TerrainNet-TA. Only temporally aggregating the stereo depths creates an incomplete elevation map due to occlusion. A non-learning-based interpolation cannot leverage the context to predict the occluded areas, such as the hill crest and the flat ground behind dense vegetation. In contrast, TerrainNet learns shape priors to predict occluded areas more accurately. TerrainNet is also less prone to the artefacts in the raw stereo depths.

3) *Inference Speed*: Inference speed is vital for fast and safe off-road driving. In Table I, TerrainNet is $7\times$ faster than LSS. The significant speedup comes from the direct projection of point features, which avoids the cost of projecting the features to every location on a depth ray. Adding stereo depth to the input has a minor impact on the speed because it only introduces a few extra convolution layers. Similarly, the temporal aggregation module has a relatively low overhead.

B. Ablation Study

1) *Architecture*: To understand the importance of each component in TerrainNet, we perform an ablation study by training a few alternative models with some of the components disabled. We summarize the results in Table II and have the following observations:

- 1) Without soft quantization (“*No soft quantization*”), the semantic mIoU degrades by 2 points, indicating the differentiable projection is crucial for correcting the end-to-end projection error from image space to BEV space. It has a small impact on elevation accuracy likely due to the terrain being naturally smooth.
- 2) The “*No z*” variant removes the elevation feature completely and only projects the image features. This results in a large error in ground elevation.
- 3) Removing the depth completion network and only using the raw stereo depth (“*No depth completion*”) degrades both semantic and elevation accuracy by a large margin.
- 4) We also trained a model using the ground-truth depth for the projection (“*Full model + GT depth*”). It has a sizeable improvement both in semantic mIoU (2 points) and elevation error.

The results show that every design decision made in TerrainNet is crucial for its performance. It also indicates that improving depth prediction accuracy is an effective way to improve the terrain model.

2) *Pseudo-labeling*: Since we train a pseudo-labeling model to generate the semantic BEV labels for the unlabeled training frames, we compare this scheme with two alternatives that do not leverage pseudo-labeling: 1) using only labeled data (about 4500 training frames), and 2) using the whole training set but not applying the semantic loss to unlabeled frames (note that elevation loss is applied to every frame).

Table III compares the pseudo-labeling scheme with the two alternatives on the validation set (we do not apply pseudo-labeling to the validation set). Pseudo-labeling provides the largest gain in semantic prediction (3.7 points) but at a small cost in elevation accuracy. Tuning the weight between the semantic and elevation loss can potentially mitigate this problem, which we leave as future work.

3) *Map size*: We choose $50\text{ m} \times 50\text{ m}$ maps for our main comparison because the cameras on the vehicle are tilted downwards such that the image content beyond 25 m gets heavily compressed and all models struggle at larger maps. Nevertheless, our model can scale to $100\text{ m} \times 100\text{ m}$ maps better than Lift-Splat-Shoot. In Table IV we compare Lift-Splat-Shoot, TerrainNet, and TerrainNet-TA on $100\text{ m} \times 100\text{ m}$ maps with RGB-D inputs.

	mIoU \uparrow	Ground Elevation (MAE) \downarrow
No soft quantization	0.449	0.254
No z	0.451	0.628
No depth completion	0.432	0.314
Full model	0.469	0.244
Full model + GT depth	0.496	0.232

TABLE II: Ablation study with TerrainNet and RGB-D inputs.

Training data	mIoU \uparrow	Ground Elevation (MAE) \downarrow
Only manually labeled	0.432	0.281
Manually labeled + All elev.	0.427	0.241
Pseudo labeled + All elev.	0.469	0.244

TABLE III: Effects of training data on the terrain modeling accuracy.

We see a similar trend in semantic accuracy: TerrainNet-TA performs the best in semantic prediction. In terms of scalability, TerrainNet is now $10\times$ faster than LSS. TerrainNet projects each pixel to one map location, so its projection procedure is independent of map size.

C. Planning with TerrainNet

To show that the output of TerrainNet can be effectively used by a planner for off-road navigation, we also ran experiments with a planner in the loop.

Planner. We use MPPI [53], a sampling-based model predictive control algorithm, as our planner as it is effective for high-speed off-road driving [52, 53]. At each time step, we perform 3000 rollouts with a kinematic bicycle model and evaluate the cost of each rollout using the terrain features along the trajectory, as described below. Then, we compute the optimized control trajectory as a weighted average of the rollouts, with the weights computed from the rollouts' aforementioned costs.

Cost function. The terrain cost C of a trajectory τ is a sum of costs evaluated at each state: $C(\tau) = \sum_{t=1}^T G(s_t) + M(s_t)$, where T is the planning horizon, s_t is the planar vehicle state (location on the map, velocity, and heading) at time step t , $G(s)$ is the distance from state s to some desired goal state, and $M(s)$ is the terrain cost at s .

The terrain cost M incorporates both semantic and geometric information based on the capabilities of the vehicle. There are several ways to do this. For the static comparative experiments in Figure 8, we compute it as follows:

$$\begin{aligned}
M(s) &= c_{\text{semantic}}^{\text{ground}} + c_{\text{semantic}}^{\text{ceiling}} + c_{\text{elevation}} \\
c_{\text{semantic}}^{\text{ground}} &= (1 + \alpha_1(h_{\text{max}} - h_{\text{min}}))\gamma_{\text{ground}}^{\text{T}} C_{\text{ground}} \\
c_{\text{semantic}}^{\text{ceiling}} &= (1 + \alpha_2(h_{\text{min}} + h_{\text{clearance}} - h_{\text{ceiling}}))\gamma_{\text{ceiling}}^{\text{T}} C_{\text{ceiling}} \\
c_{\text{elevation}} &= \beta_1\theta_{\text{roll}}^2 + \beta_2\theta_{\text{pitch}}^2,
\end{aligned}$$

where α_1 , α_2 , β_1 , and β_2 are scalar tuning parameters, γ_{ground} and γ_{ceiling} are tuning vectors of costs for each semantic class, and θ_{roll} and θ_{pitch} are the estimated roll and pitch angles of the vehicle based on the state s and the elevation map from TerrainNet. The height multiplier for semantic costs helps account for the fact that within any given semantic class, larger obstacles tend to be less traversable. The elevation cost penalizes large pitch and roll angles. Large pitch angles may be too steep to ascend, whereas large roll angles may cause the vehicle to tip over. These angles depend on vehicle heading

Method	mIoU \uparrow	Ground Elevation (MAE) \downarrow	Time (ms) \downarrow
LSS	0.399	0.628	402
TerrainNet	0.419	0.586	39
TerrainNet-TA	0.464	0.622	42

TABLE IV: Results on $100\text{ m} \times 100\text{ m}$ maps with RGB-D inputs.

Method	Avg. HD \downarrow	Avg. CD \downarrow	Time (ms) \downarrow
SimpleBEV	2.837	0.404	113
LSS	2.760	0.387	207
TerrainNet	2.658	0.371	28
TerrainNet-TA	2.758	0.384	29

TABLE V: Planning evaluation with open-loop MPPI. Inputs are RGB-D. We compute the Average symmetric Hausdorff Distance and the Average Cost Difference between the optimized control trajectories on the predicted and the ground-truth costmaps.

and can be computed using h_{min} (or h_{max}) at the locations of the wheels.

For the real-world experiments (Figure 9), the terrain cost incorporates more physics:

$$M(s) = v^2 c_{\text{semantic}}^{\text{ground}} + C_{\text{lethal}} \mathbf{1}[c_{\text{rollover}} > \delta] \quad (8)$$

$$c_{\text{rollover}} = |\kappa v^2 + g \sin(\theta_{\text{roll}})| / \cos(\theta_{\text{roll}}), \quad (9)$$

where C_{lethal} is the lethal cost, $\mathbf{1}[\cdot]$ is the indicator function, v is vehicle velocity, κ is trajectory curvature, g is the acceleration due to gravity, and δ is a tuning parameter depending on vehicle geometry. Multiplying the semantic cost $c_{\text{semantic}}^{\text{ground}}$ by v^2 reflects the fact that collisions are more dangerous at higher speeds. c_{rollover} is the risk of rolling the vehicle based on speed, ground slope, and how sharply the vehicle is turning.

Static planning comparison. First, we set up an experiment on the validation dataset by running the MPPI planner on the costmap computed from the output of TerrainNet and other baselines (all with stereo inputs). The navigation task is as follows: we choose 12 waypoints in front of the vehicle at a distance of 25 m, spanning from -60° to 60° . Then, we run MPPI to plan a trajectory to each waypoint. Table V summarizes the results, and Figure 8 visualizes example costmaps and trajectories. We observe the same trend where TerrainNet with RGB-D inputs performs the best. One interesting observation is that TerrainNet-TA performs worse than TerrainNet. This is likely caused by the planner being more sensitive to certain terrain features than others, and that the vehicle mostly driving forward (making history less useful). Currently, for generality, we do not discriminate between different terrain classes, but it would be beneficial to weigh their influence on planning during training. We leave this as future work.

Real-world navigation. TerrainNet has been integrated with an off-road autonomous driving stack and tested in real-world environments. Limited by the season, we tested TerrainNet in an off-road environment covered with deep snow and steep slopes. Since our dataset does not contain snow, we annotated 15 LiDAR scans with snow and finetuned our model with an updated ontology (including *snow* and *snow trail*). Figure 9 (best viewed zoomed in) shows sample costmaps, elevation maps, and planner output from an autonomous run. The system completed a 1.1 km run with a maximum speed of 7 m/s (average of 3.2 m/s) with two human interventions. The interventions were mainly due to errors in odometry from

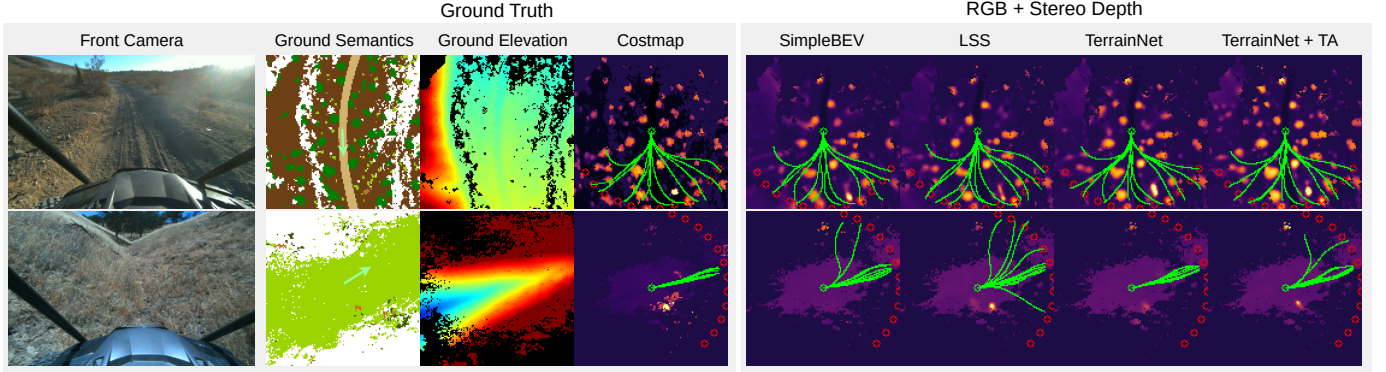


Fig. 8: Visualization of MPPI control trajectories on two different terrains. Green arrow indicates vehicle’s heading. Brighter areas correspond to higher costs. All models use RGB with stereo depth. **Top:** On-trail driving. The trajectories may go through non-lethal vegetation but do avoid the tree on the front left of the vehicle. **Bottom:** Driving in a steep valley. Due to the steep slopes, the trajectories should stay in the valley. We do not visualize the elevation cost here due to its dependence on vehicle heading.

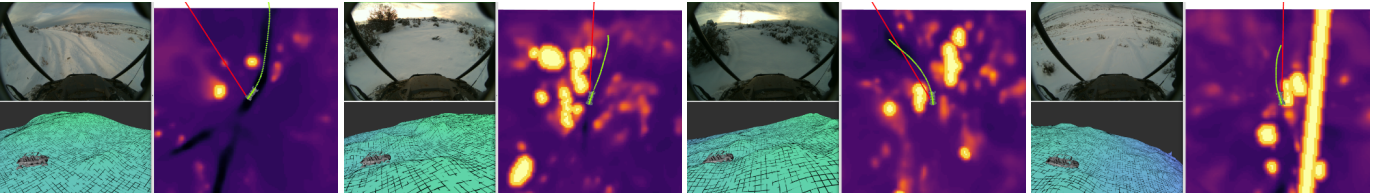


Fig. 9: Qualitative results from a closed-loop real-world experiment. Each figure shows the front RGB camera (non-rectified), elevation map (showing h_{\max}), and semantic costmap (showing c_{semantic}). Overlaid on the semantic costmap, the red line shows the direction of the goal point, and the green line shows the local plan. From the left: **(1) Driving on a trail with a hill to the left side.** Note that the trail is lower cost, so the local planner moves faster (as shown by the long green line). **(2) Driving off trail up a hill.** Note that the tall bushes on the left are very high cost. **(3) Driving uphill with many wheel slips on the snow.** The costmap and elevation map are “smeared” along the direction of travel, which highlights a limitation of temporal aggregation when odometry is poor. **(4) Approaching a steep downhill slope.** Although the slope itself is not yet directly visible, TerrainNet predicts a steep elevation map. The local planner slows down to reduce the risk of rollover. (The high-cost straight line through the costmap is a manually GPS-specified keep-out zone to mark the location of a nearby fence.)

wheel slips. TerrainNet ran at 20 Hz onboard the vehicle. Due to additional overhead from the system, TerrainNet provided map updates at 10 Hz. See the website for more details.

VII. DISCUSSION

LiDAR point density. One question a reader may ask is whether increasing the number of LiDARs is sufficient for building a dense map at high speed. While more LiDARs increase the point density, it also increases the cost of the system. Compared to typical deployment scenarios where robots may only have one LiDAR, cameras provide much higher pixel density at a lower cost. Cameras also provide additional advantages such as higher reliability, stealth, and being less affected by small particles such as dust. Nonetheless, it would be interesting to compare LiDAR-based mapping and camera-based mapping in challenging environments.

Limitations. While TerrainNet predicts high-fidelity terrain semantics and elevations, it has several limitations to be addressed:

- It may not generalize well to an environment that is substantially different from the training environments. This is common in learning-based systems. To improve the generalization, we can increase the diversity of the datasets and apply domain adaptation techniques.
- It may miss fine-grained features such as small, lethal rocks. Due to hardware limitations, the cameras are

not precisely synchronized with the LiDARs (note that LiDARs are used to generate training data). Hence, there exists a small but noisy misalignment between the projected camera features and the ground-truth maps. This makes it harder for the model to learn fine-grained details of the terrain. To address this, we can improve the sensor alignment and make the model less sensitive to the localization error of small objects.

- TerrainNet does not predict uncertainty. This can lead to dangerous behaviors, such as the model predicting the other side of a cliff as a ramp. We can incorporate uncertainty estimation or risk assessment as done by Cai et al. [8] and Fan et al. [15] to make the model risk-aware.

VIII. CONCLUSION

In this paper, we designed and implemented TerrainNet, the first off-road, camera-only perception system for joint BEV semantic and geometric terrain mapping. We demonstrated TerrainNet’s accuracy and efficiency in terrain perception and successfully deployed the model as part of a navigation system for a variety of challenging off-road environments. In future work, we aim to include uncertainty estimates in the output, learn costmaps directly from expert demonstrations to remove manually engineered cost parameters, and more effectively leverage unlabeled data for learning.

ACKNOWLEDGMENTS

This research was developed with funding from the Defense Advanced Research Projects Agency (DARPA), ARL SARA CRA W911NF-20-2-0095, and ARL SARA CRA W911NF-21-2-0190.

REFERENCES

- [1] Multisense stereo cameras. URL <https://www.carnegierobotics.com/products/multisense-s27>.
- [2] Polaris RZR. URL <https://rzt.polaris.com>.
- [3] PyTorch Smooth-L1 Loss. URL <https://pytorch.org/docs/stable/generated/torch.nn.SmoothL1Loss.html>.
- [4] Max Bajracharya, Andrew Howard, Larry H Matthies, Benyang Tang, and Michael Turmon. Autonomous off-road navigation with end-to-end learning for the LAGR program. *Journal of Field Robotics*, 26(1):3–25, 2009.
- [5] Max Bajracharya, Jeremy Ma, Matt Malchano, Alex Perkins, Alfred A Rizzi, and Larry Matthies. High fidelity day/night stereo mapping with vegetation and negative obstacle detection for vision-in-the-loop walking. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 3663–3670. IEEE, 2013.
- [6] Marcelo Bertalmio, Andrea L Bertozzi, and Guillermo Sapiro. Navier-Stokes, fluid dynamics, and image and video inpainting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2001.
- [7] Paulo V. K. Borges, Thierry Peynot, Sisi Liang, Bilal Arain, Matthew Wildie, Melih G. Minareci, Serge Lichman, Garima Samvedi, Inkyu Sa, Nicolas Hudson, Michael Milford, Peyman Moghadam, and Peter Corke. A survey on terrain traversability analysis for autonomous ground vehicles: Methods, sensors, and challenges. *Field Robotics*, 2(1):1567–1627, 2022.
- [8] Xiaoyi Cai, Michael Everett, Jonathan Fink, and Jonathan P How. Risk-aware off-road navigation via a learned speed distribution map. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 2931–2937. IEEE, 2022.
- [9] Xiaoyi Cai, Michael Everett, Lakshay Sharma, Philip R Osteen, and Jonathan P How. Probabilistic traversability model for risk-aware motion planning in off-road environments. *arXiv preprint arXiv:2210.00153*, 2022.
- [10] Zhe Chen, Yuchen Duan, Wenhai Wang, Junjun He, Tong Lu, Jifeng Dai, and Yu Qiao. ViT-Adapter: Exploring plain vision transformer for accurate dense predictions. In *International Conference on Learning Representations*, 2023.
- [11] Kashyap Chitta, Aditya Prakash, and Andreas Geiger. NEAT: Neural attention fields for end-to-end autonomous driving. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15793–15803, 2021.
- [12] Nitish Dashora, Daniel Shin, Dhruv Shah, Henry Leopold, David Fan, Ali Agha-Mohammadi, Nicholas Rhinehart, and Sergey Levine. Hybrid imitative planning with geometric and predictive costs in off-road environments. In *International Conference on Robotics and Automation (ICRA)*, pages 4452–4458. IEEE, 2022.
- [13] Tom Duckett, Simon Pearson, Simon Blackmore, Bruce Grieve, and Melvyn Smith. Agricultural robotics: The future of robotic agriculture, 2018. URL <https://uwe-repository.worktribe.com/output/866226>.
- [14] Alberto Elfes. Using occupancy grids for mobile robot perception and navigation. *Computer*, 22(6):46–57, 1989.
- [15] David D Fan, Ali-Akbar Agha-Mohammadi, and Evangelos A Theodorou. Learning risk-aware costmaps for traversability in challenging environments. *IEEE Robotics and Automation Letters*, 7(1):279–286, 2021.
- [16] Péter Fankhauser, Michael Bloesch, and Marco Hutter. Probabilistic terrain mapping for mobile robots with uncertain localization. *IEEE Robotics and Automation Letters*, 3(4):3019–3026, 2018.
- [17] Bianca Forkel and Hans-Joachim Wuensche. Dynamic resolution terrain estimation for autonomous (dirt) road driving fusing lidar and vision. In *Intelligent Vehicles Symposium (IV)*, pages 1181–1187. IEEE, 2022.
- [18] Bianca Forkel, Jan Kallwies, and Hans-Joachim Wuensche. Probabilistic terrain estimation for autonomous off-road driving. In *International Conference on Robotics and Automation (ICRA)*, pages 13864–13870. IEEE, 2021.
- [19] Tianrui Guan, Zhenpeng He, Ruitao Song, Dinesh Manocha, and Liangjun Zhang. TNS: Terrain traversability mapping and navigation system for autonomous excavators. In *Robotics: Science and Systems*, 2022.
- [20] Yutao Han, Jacopo Banfi, and Mark Campbell. Planning paths through unknown space by imagining what lies therein. In *Conference on Robot Learning*, 2021.
- [21] Adam W Harley, Zhaoyuan Fang, Jie Li, Rares Ambrus, and Katerina Fragkiadaki. Simple-BEV: What really matters for multi-sensor BEV perception? In *International Conference on Robotics and Automation (ICRA)*. IEEE, 2023.
- [22] Noureldin Hendy, Cooper Sloan, Feng Tian, Pengfei Duan, Nick Charchut, Yuesong Xie, Chuang Wang, and James Philbin. FISHING Net: Future inference of semantic heatmaps in grids. *arXiv preprint arXiv:2006.09917*, 2020.
- [23] Wolfgang Hess, Damon Kohler, Holger Rapp, and Daniel Andor. Real-time loop closure in 2D LIDAR SLAM. In *International Conference on Robotics and Automation (ICRA)*, pages 1271–1278. IEEE, 2016.
- [24] Maximilian Jaritz, Raoul De Charette, Emilie Wirbel, Xavier Perrotton, and Fawzi Nashashibi. Sparse and dense data with CNNs: Depth completion and semantic segmentation. In *International Conference on 3D Vision (3DV)*, pages 52–60. IEEE, 2018.
- [25] Peng Jiang, Philip Osteen, Maggie Wigness, and Srikanth Saripalli. RELLIS-3D dataset: Data, benchmarks and analysis. In *International Conference on Robotics and Automation (ICRA)*. IEEE, 2021.
- [26] Alonzo Kelly, Anthony Stentz, Omead Amidi, Mike Bode, David Bradley, Antonio Diaz-Calderon, Mike Hoppold, Herman Herman, Robert Mandelbaum, Tom Pilarski, Pete

- Rander, Anthony Stenz, Nick Vallidis, and Randy Warner. Toward reliable off road autonomous vehicles operating in challenging environments. *The International Journal of Robotics Research*, 25(5-6):449–483, 2006.
- [27] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015.
- [28] Eric Krotkov and John Blitch. The Defense Advanced Research Projects Agency (DARPA) tactical mobile robotics program. *The International Journal of Robotics Research*, 18(7):769–776, 1999.
- [29] Alex H Lang, Sourabh Vora, Holger Caesar, Lubing Zhou, Jiong Yang, and Oscar Beijbom. PointPillars: Fast encoders for object detection from point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019.
- [30] Dong-Hyun Lee. Pseudo-Label: The simple and efficient semi-supervised learning method for deep neural networks. In *Workshop on Challenges in Representation Learning, ICML*, volume 3, page 896, 2013.
- [31] Zhiqi Li, Wenhai Wang, Hongyang Li, Enze Xie, Chonghao Sima, Tong Lu, Yu Qiao, and Jifeng Dai. BEVFormer: Learning bird’s-eye-view representation from multi-camera images via spatiotemporal transformers. In *European Conference on Computer Vision*, 2022.
- [32] Roberto Manduchi, Andres Castano, Ashit Talukder, and Larry Matthies. Obstacle detection and terrain classification for autonomous off-road navigation. *Autonomous Robots*, 18(1):81–102, 2005.
- [33] Daniel Maturana, Po-Wei Chou, Masashi Uenoyama, and Sebastian Scherer. Real-time semantic mapping for autonomous off-road navigation. In *Field and Service Robotics: Results of the 11th International Conference*, pages 335–350. Springer, 2018.
- [34] Takahiro Miki, Lorenz Wellhausen, Ruben Grandia, Fabian Jenelten, Timon Homberger, and Marco Hutter. Elevation mapping for locomotion and navigation using GPU. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 2273–2280. IEEE, 2022.
- [35] Jonah Philion and Sanja Fidler. Lift, splat, shoot: Encoding images from arbitrary camera rigs by implicitly unprojecting to 3D. In *European Conference on Computer Vision*, 2020.
- [36] Rui Qian, Divyansh Garg, Yan Wang, Yurong You, Serge Belongie, Bharath Hariharan, Mark Campbell, Kilian Q Weinberger, and Wei-Lun Chao. End-to-end pseudo-LiDAR for image-based 3D object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5881–5890, 2020.
- [37] Cody Reading, Ali Harakeh, Julia Chae, and Steven L Waslander. Categorical depth distribution network for monocular 3D object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8555–8564, 2021.
- [38] Thomas Roddick and Roberto Cipolla. Predicting semantic map representations from images using pyramid occupancy networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11138–11147, 2020.
- [39] Abbas Sadat, Sergio Casas, Mengye Ren, Xinyu Wu, Pranaab Dhawan, and Raquel Urtasun. Perceive, predict, and plan: Safe motion planning through interpretable semantic representations. In *European Conference on Computer Vision*, 2020.
- [40] Avishkar Saha, Oscar Mendez, Chris Russell, and Richard Bowden. Translating images into maps. In *2022 International Conference on Robotics and Automation (ICRA)*, pages 9200–9206. IEEE, 2022.
- [41] Vojtěch Šalanský, Karel Zimmermann, Tomáš Petříček, and Tomáš Svoboda. Pose consistency KKT-loss for weakly supervised learning of robot-terrain interaction model. *IEEE Robotics and Automation Letters*, 6(3): 5477–5484, 2021.
- [42] Fabian Schilling, Xi Chen, John Folkesson, and Patric Jensfelt. Geometric and visual terrain classification for autonomous mobile navigation. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2017.
- [43] Christos Sevastopoulos and Stasinios Konstantopoulos. A survey of traversability estimation for mobile robots. *IEEE Access*, 10:96331–96347, 2022.
- [44] Amirreza Shaban, Xiangyun Meng, JoonHo Lee, Byron Boots, and Dieter Fox. Semantic terrain classification for off-road autonomous driving. In *Conference on Robot Learning*, pages 619–629. PMLR, 2022.
- [45] David Silver, J Andrew Bagnell, and Anthony Stentz. Learning from demonstration for autonomous navigation in complex unstructured terrain. *The International Journal of Robotics Research*, 29(12):1565–1592, 2010.
- [46] Maximilian Stölzle, Takahiro Miki, Levin Gerdes, Martin Azkarate, and Marco Hutter. Reconstructing occluded elevation information in terrain maps with self-supervised learning. *IEEE Robotics and Automation Letters*, 7(2): 1697–1704, 2022.
- [47] Tarlan Suleymanov, Lina Maria Paz, Pedro Piniés, Geoff Hester, and Paul Newman. The path less taken: A fast variational approach for scene segmentation used for closed loop control. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 3620–3626. IEEE, 2016.
- [48] Andrew Thoesen and Hamid Marvi. Planetary surface mobility and exploration: A review. *Current Robotics Reports*, 2(3):239–249, 2021.
- [49] Rudolph Triebel, Patrick Pfaff, and Wolfram Burgard. Multi-level surface maps for outdoor terrain mapping and loop closing. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 2276–2282. IEEE, 2006.
- [50] Rodrigo Ventura and Pedro U Lima. Search and rescue robots: The civil protection teams of the future. In *Third International Conference on Emerging Security Technologies*, pages 12–19. IEEE, 2012.
- [51] Brian H Wilcox. Robotic vehicles for planetary explo-

- ration. *Applied Intelligence*, 2:181–193, 1992.
- [52] Grady Williams, Paul Drews, Brian Goldfain, James M Rehg, and Evangelos A Theodorou. Aggressive driving with model predictive path integral control. In *International Conference on Robotics and Automation (ICRA)*, pages 1433–1440. IEEE, 2016.
- [53] Grady Williams, Nolan Wagener, Brian Goldfain, Paul Drews, James M Rehg, Byron Boots, and Evangelos A Theodorou. Information theoretic MPC for model-based reinforcement learning. In *International Conference on Robotics and Automation (ICRA)*, pages 1714–1721. IEEE, 2017.
- [54] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. SegFormer: Simple and efficient design for semantic segmentation with transformers. In *Advances in Neural Information Processing Systems 34*, pages 12077–12090, 2021.
- [55] Weixiang Yang, Qi Li, Wenxi Liu, Yuanlong Yu, Yuexin Ma, Shengfeng He, and Jia Pan. Projecting your view attentively: Monocular road scene layout estimation via cross-view transformation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021.
- [56] Brady Zhou and Philipp Krähenbühl. Cross-view transformers for real-time map-view semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13760–13769, 2022.
- [57] Xinge Zhu, Hui Zhou, Tai Wang, Fangzhou Hong, Yuexin Ma, Wei Li, Hongsheng Li, and Dahua Lin. Cylindrical and asymmetrical 3D convolution networks for LiDAR segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9939–9948, 2021.

APPENDIX

A. Computing Terrain Representation

This section provides more details about how we build the terrain representation.

Point cloud aggregation. Given the raw LiDAR packets and the vehicle’s ego motion, we convert the LiDAR packets into LiDAR points with the rolling shutter effects corrected. Then, we run Google Cartographer [23] on the LiDAR points (with IMU information) to compute the gravity-aligned pose for each LiDAR scan. We aggregate the LiDAR scans using their poses to get an aggregated point cloud.

Outlier removal. LiDAR may produce false returns due to the presence of dust, snowflakes, and water. These false returns manifest as noise in the aggregated point cloud. We identify these outliers in two ways:

- 1) Voxel filter. We voxelize the point cloud with a voxel size of 0.3m. We count the number of points in each voxel. If the number of points is smaller than a threshold (set to 5), we remove all the points in this voxel.
- 2) Semantic filter. When we label the LiDAR points, we ask the labelers to label the outliers as an additional *outlier*

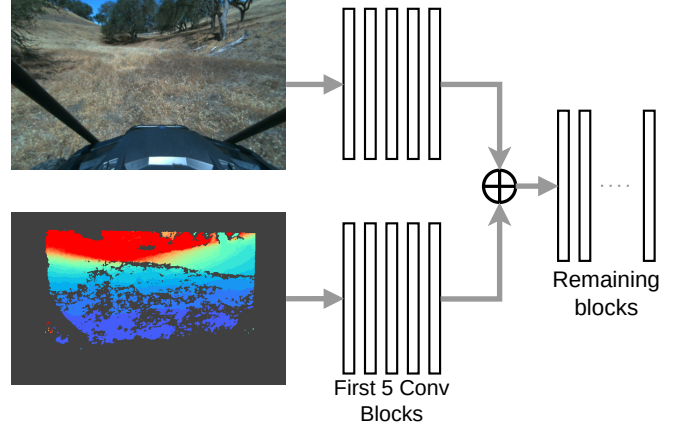


Fig. 10: The image backbone uses separate branches to process the RGB and depth inputs, respectively. Then, the outputs of the two branches are summed and passed to the remaining blocks in the backbone.

class. Our LiDAR segmentation network is trained with the additional outlier class. After we predict the pseudo labels for each point, we remove points classified as outliers.

B. Network Architecture

Backbone. Our image backbone follows that of Lift-Splat-Shoot[35]. There are two main differences:

- 1) For RGB-D inputs, we duplicate the first five convolutional blocks to process the RGB and depth images separately (Figure 10). This has a noticeable improvement over passing a 4-channel RGB-D image directly to the backbone.
- 2) We do $8\times$ downsampling of the feature map and depth map (the original paper performs $16\times$ downsampling). This provides a good trade-off between speed and accuracy. Note that all the baselines are trained with the same image backbone with the same $8\times$ downsampling ratio.

Our input image size is 512×320 . Hence, the image backbone produces a 472-channel feature map with a spatial dimension of 64×40 .

Depth completion network. The output of the image backbone is passed to a single convolution layer with 472 input channels and D output channels with kernel size 3. D is the number of desirable depth bins. For $50 \text{ m} \times 50 \text{ m}$ maps, $D = 128$. For $100 \text{ m} \times 100 \text{ m}$ maps, $D = 256$.

Terrain embedding network. For each back-projected point (x, y, z) , we first apply an MLP with layer sizes $(1, 64, 32)$ to z to get the 32-dimensional elevation embedding f_{elev} . Then, we concatenate f_{elev} with the 472-channel image feature f_{sem} and apply another MLP with layer sizes $(504, 96)$ to get the 96-dimensional terrain embedding f . ReLU is applied after each MLP layer.

Temporal aggregation layer. The temporal aggregation layer is a single ConvGRU layer with 96 input channels, 96 output channels, and a kernel size of 1. To train TerrainNet-TA, we first train the non-recurrent version, and then insert the

recurrent layer. We freeze the weights of the model except the recurrent layer.

Inpainting network. The inpainting net follows the architecture of Lift-Splat-Shoot’s BEV encoder network. The main difference is that we create a separate decoding head consisting of an upsampling layer and two convolution layers to predict a specific terrain layer.

C. Data Annotation

We create our data annotation tool (Figure 11) to obtain ground truth semantically labeled point clouds. We aggregate 30 to 50 LiDAR scans for each selected frame to get a dense point cloud for labeling. We provide camera images for every

LiDAR scan for reference so that a labeler can cross-reference to verify that the labels are correct.

D. Additional dataset examples

Figure 12 shows additional dataset examples to highlight the diversity of the datasets. We collected our datasets from 3 distinctive geographic locations across the year.

E. More qualitative results

Figure 13 shows other example predictions of TerrainNet and other baselines on the validation sets. Figure 14 visualizes different predicted 3D terrains of TerrainNet-TA on the validation sets. For more details, including videos, please visit the website.

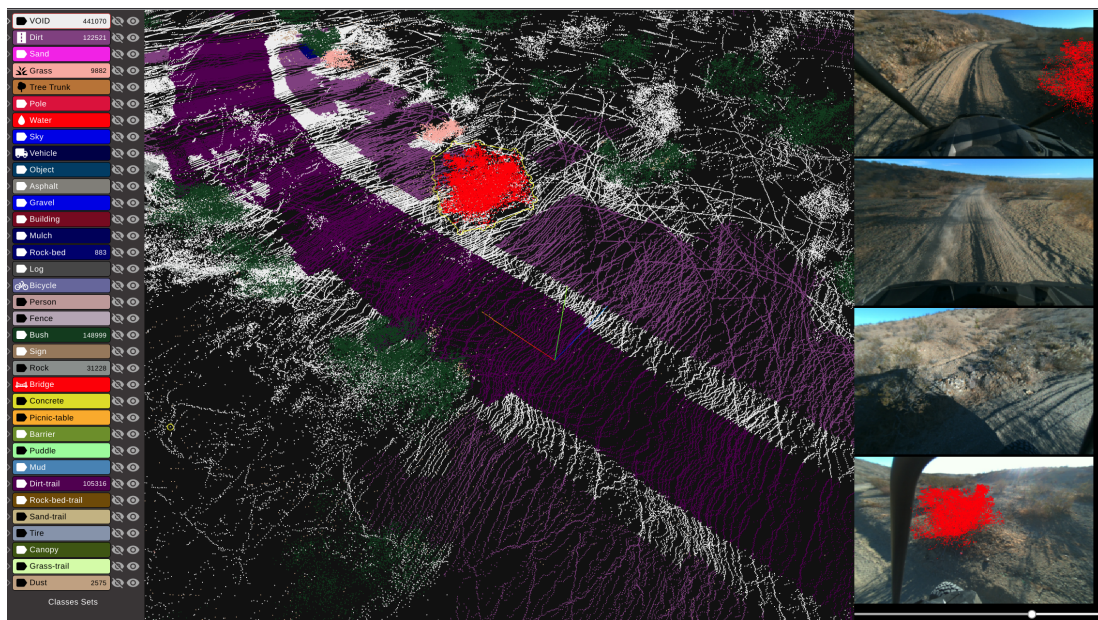


Fig. 11: **Annotation tool.** Labelers can select points in the point cloud and see their corresponding image projections. Here we show a bush is highlighted both in the point cloud and in the bottom image. This helps the labelers to verify the semantic class of those points. Moreover, labelers only annotate points that they are confident about. We leave points that are difficult to identify unlabeled.

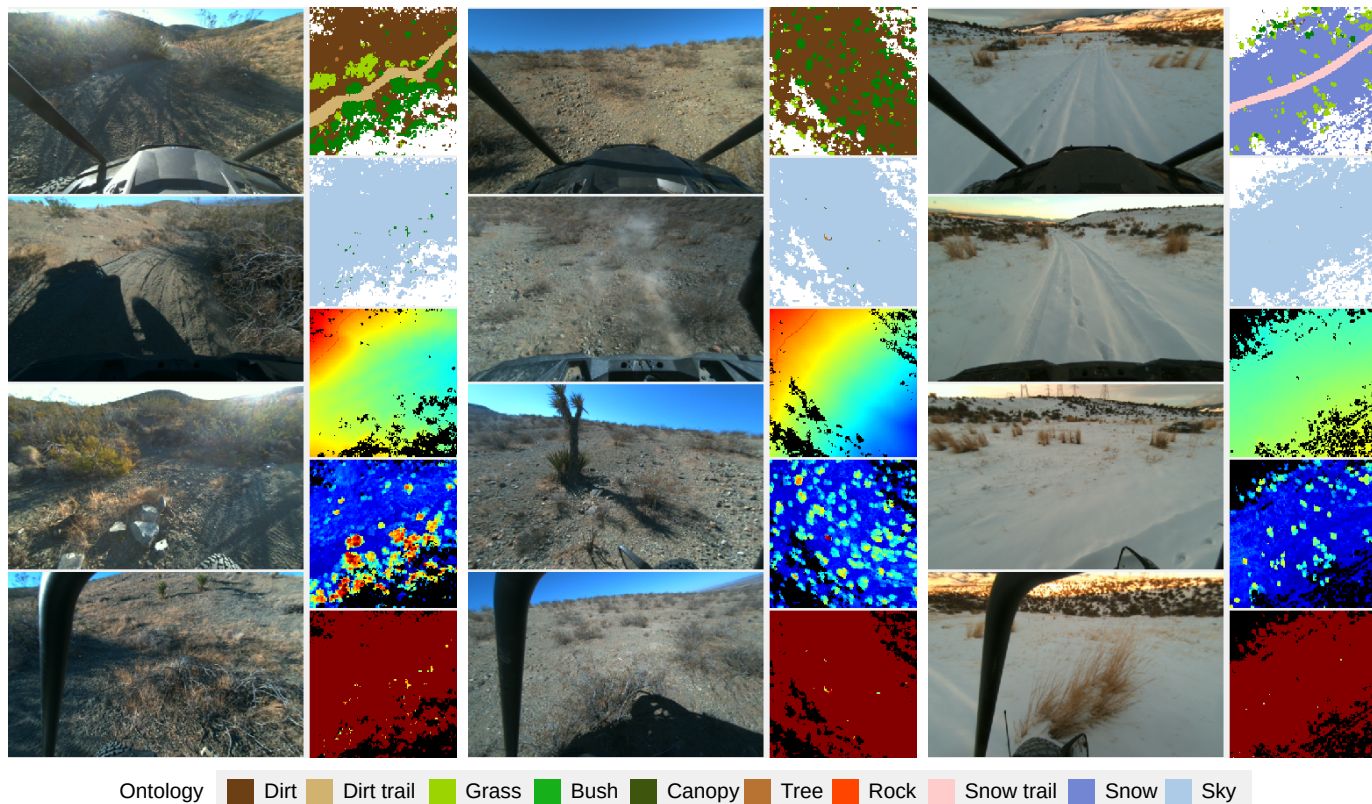


Fig. 12: Additional dataset examples. These examples were collected at different locations and in different seasons. For the real-world experiments, we finetuned the model with annotated snow data (rightmost column).

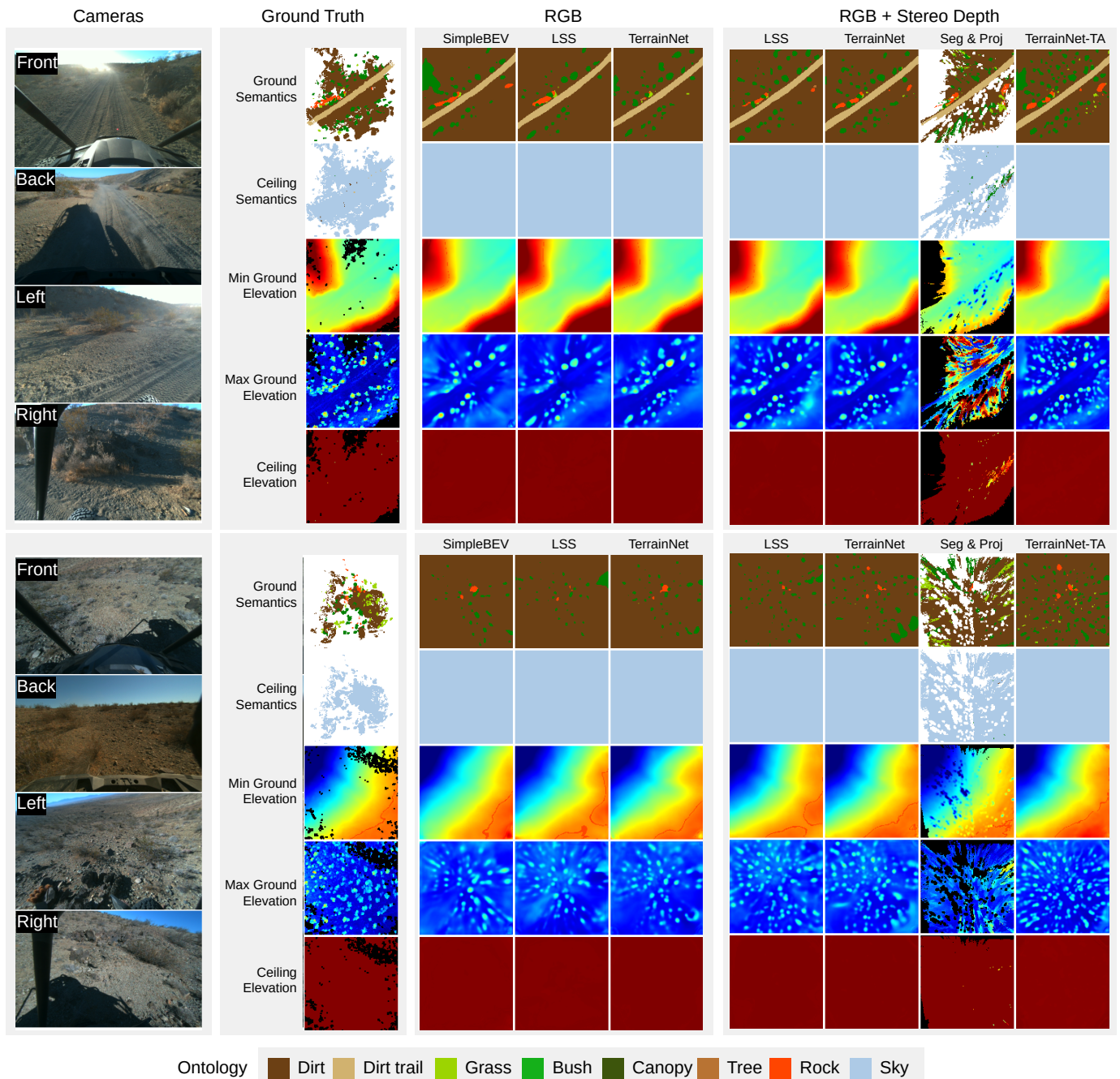


Fig. 13: Comparing the predictions of different models. **Top:** on-trail run with rocks and bushes on the two sides of the trail. **Bottom:** off-trail run with scattered rocks and bushes. These are open terrains, so the ceiling semantics consists almost entirely of *sky*.

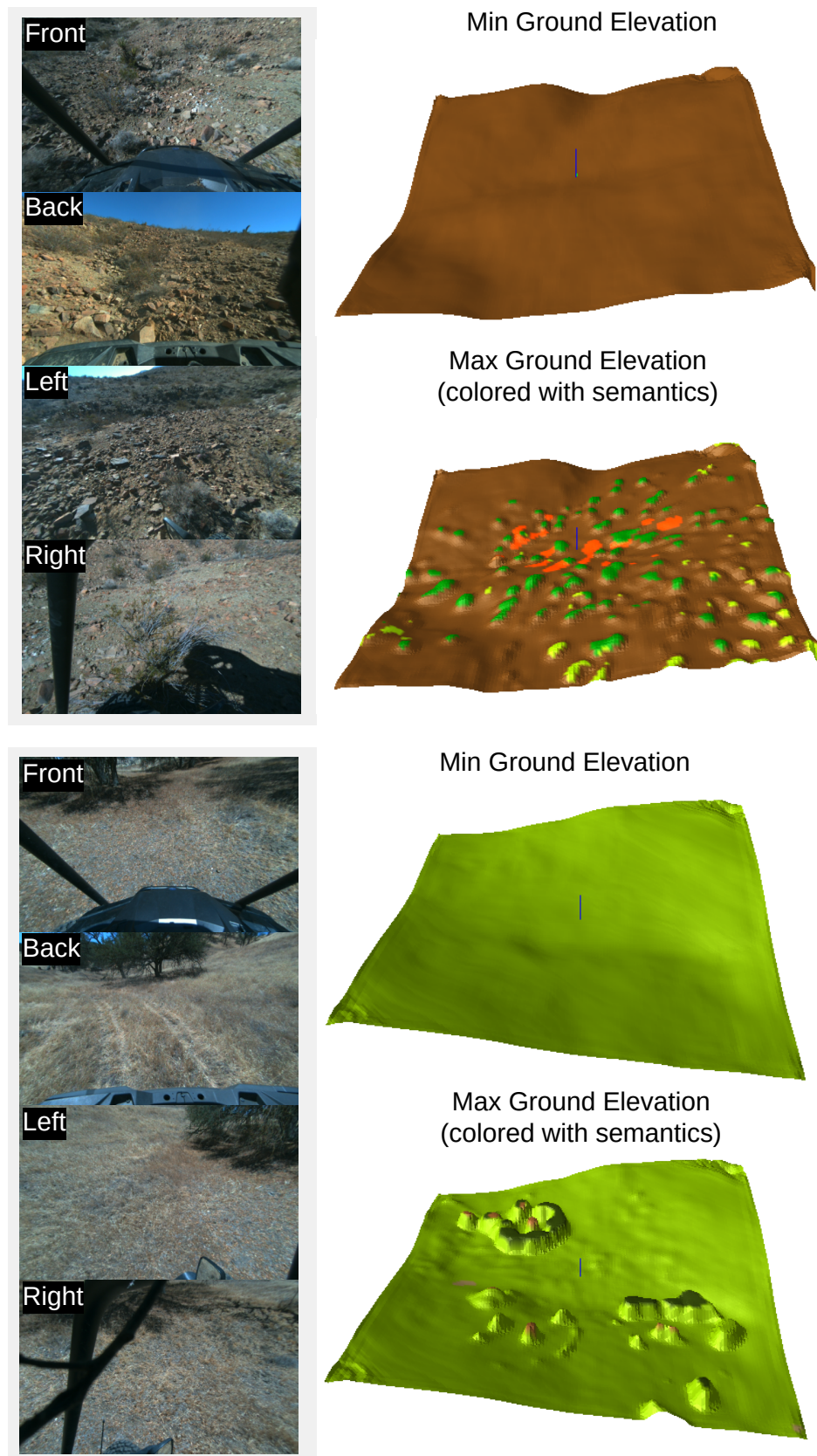


Fig. 14: **Visualizations of the predicted 3D terrain model from TerrainNet-TA.** We only show the min and max ground elevation here for clarity. Left: uneven terrain with scattered grass, bushes, and rocks. Right: a grass-covered valley with tall trees and low-hanging canopies.