

ConvBKI: Real-Time Probabilistic Semantic Mapping Network with Quantifiable Uncertainty

Joey Wilson, Yuewei Fu, Joshua Friesen, Parker Ewen
Andrew Capodiecici, Paramsothy Jayakumar, Kira Barton, and Maani Ghaffari

Abstract—In this paper, we develop a modular neural network for real-time (> 10 Hz) semantic mapping in uncertain environments, which explicitly updates per-voxel probabilistic distributions within a neural network layer. Our approach combines the reliability of classical probabilistic algorithms with the performance and efficiency of modern neural networks. Although robotic perception is often divided between modern differentiable methods and classical explicit methods, a union of both is necessary for real-time and trustworthy performance. We introduce a novel Convolutional Bayesian Kernel Inference (ConvBKI) layer which incorporates semantic segmentation predictions online into a 3D map through a depthwise convolution layer by leveraging conjugate priors. We compare ConvBKI against state-of-the-art deep learning approaches and probabilistic algorithms for mapping to evaluate reliability and performance. We also create a Robot Operating System (ROS) package of ConvBKI and test it on real-world perceptually challenging off-road driving data.

I. INTRODUCTION

Robotic perception is at a crossroads between classical, probabilistic methods and modern, implicit deep neural networks. Classical probabilistic solutions offer reliable and trustworthy performance at the expense of longer run-times and unoptimized, hand-tuned parameters. In contrast, modern deep learning methods achieve higher performance and improved latency due to optimization in an implicit space but lose the ability to generalize to new data. However, probabilistic and deep learning methods are not intrinsically opposed and may be combined through methodical design.

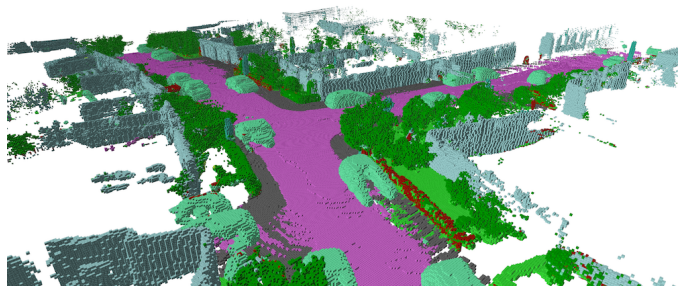
Perception is the general study of understanding sensory information, and mapping is an important subset which seeks to consolidate sensory input into an informative world model. Maps are capable of representing high levels of scene understanding by storing multiple modalities of information, including semantic labels, dynamic motion, and topological relations. Although some works propose to discard a world model in favor of mapless end-to-end control [1], maps are still widely used due to their reliability and interpretable nature with shared human-robot scene understanding. Additionally, mapping within a robotic framework enables a

DISTRIBUTION STATEMENT A. Approved for public release; distribution is unlimited. OPSEC# 7836

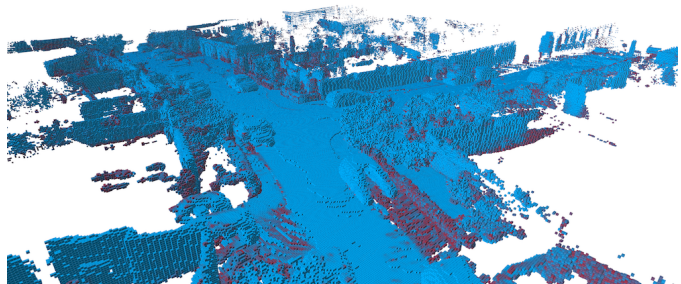
J. Wilson, Y. Fu, J. Friesen, P. Ewen, K. Barton, and M. Ghaffari are with the University of Michigan, Ann Arbor, MI 48109, USA. {wilsonjv, ywfu, friesej}@umich.edu, {pewen, bartonkl, maanigj}@umich.edu

A. Capodiecici is with Neya Systems Division, Applied Research Associates, Warrendale, PA 15086, USA. acapodiecici@neyarobotics.com

P. Jayakumar is with the US Army DEVCOM Ground Vehicle Systems Center, Warren, MI 48397, USA. paramsothy.jayakumar.civ@army.mil



(a) Semantic Categories



(b) Variance

Fig. 1: Example output of our method (ConvBKI) on Semantic KITTI sequence 0 [2]. ConvBKI recurrently updates and maintains a semantic map with uncertainty in real-time by reframing the Bayesian update step as a depthwise separable convolution operation. Voxels encode the expected semantic category and variance as concentration parameters of the Dirichlet distribution, which is the conjugate prior of the categorical distribution and multinomial distribution. The input to the network is a 3D point cloud from a stereo-camera or LiDAR sensor. The points are labeled with semantic classification probabilities by a semantic segmentation network which provides per-point predictions. ConvBKI updates the semantic belief of each voxel by spatially weighting the input semantic point clouds through a unique geometric kernel learned for each semantic category. Fig. 1a shows the maximum likelihood semantic prediction per voxel, while Fig. 1b shows the variance per voxel calculated with Eq. (10). Variance is linearly translated to RGB colors, where blue indicates voxels with a low variance and red indicates a high variance.

modular schema for perception, planning, and control which avoids formulating overly complex and ill-posed problems.

One common mapping strategy is occupancy mapping which uses 3-dimensional cubic voxels in a grid-like structure to describe the environment with binary labels. While initial maps only contained binary free or occupied labels [3], later works attempted to expand upon the scene understanding of mapping algorithms by incorporating semantic labels obtained by state-of-the-art semantic segmentation neural networks [4], [5]. Semantic labels within the map can be used by downstream planners to complete tasks such as staying on a road or lifting a cup, instead of merely avoiding occupied space. Other work has also sought to improve the rigid, discrete structure

of voxels with continuous geometry by surfels [4], [6], signed distance functions [7], meshes [8], [9] and more, although these approaches can often be represented using occupancy maps [10].

In real-world operating conditions, robots must contend with uncertain inputs corrupted by noisy sensors and limited observations. In applications where safety is paramount, such as autonomous driving, it is vital to have a measure of certainty as well as algorithms robust to new, unseen data. For example, when driving off-road, mobile robots must contend with negative obstacles in addition to positive obstacles. Whereas positive obstacles are the presence of occupied space, a negative obstacle presents a unique challenge due to the lack of a driveable surface, such as a hole or cliff, which the mobile robot cannot traverse. Without uncertainty quantification, a mobile robot may incorrectly assert the absence of an obstacle or the presence of a road instead of slowing down and investigating before proceeding.

Early approaches to probabilistic robotic mapping were based on the incorporation of input through Bayesian inference methods such as Gaussian Process Occupancy Maps [11]. However, Gaussian Processes have a cubic time complexity with respect to training samples due to expensive matrix inversion. Later Bayesian methods sought to improve the inference speed through methods such as Bayesian Generalized Kernel Inference [12], which approximate distributions at model selection [13] through kernel functions. However, real-time operation is still a challenge for many hand-crafted probabilistic mapping approaches.

In contrast, recent mapping approaches have established end-to-end world modeling architectures where high levels of scene understanding are stored in a hidden state and updated recurrently across time [14]–[16]. While implicit maps are able to achieve high performance with limited latency due to efficient and parallel modeling, they still encounter training difficulties and challenges when exposed to new data outside of their training set. Additionally, supervised approaches face the challenge of obtaining data sets which can be particularly challenging for dynamic mapping [17], [18].

In this paper, we combine probabilistic and differentiable mapping approaches to obtain a real-time (> 10 Hz) 3D semantic mapping network which balances the efficiency of deep learning approaches with the reliability of classical methods by operating on explicit probability distributions in a modular, parallelized architecture. We demonstrate that ConvBKI is able to perform comparably with deep learning approaches on data sets it has been trained on and transfer more readily to unseen data sets, including a new, challenging off-road driving data set.

A. Contributions

We pose the Bayesian semantic mapping problem as a differentiable neural network layer (ConvBKI) in order to achieve the accelerated inference rates and optimized performance of deep learning while maintaining the ability to quantify uncertainty in closed form. Compared to previous semantic mapping works such as Semantic Fusion [4] and Semantic

BKI [5], ConvBKI also performs semantic mapping in a dense voxel space with Bayesian methods but reformulates the map inference as an end-to-end neural network with differentiable parameters. We demonstrate the improved performance and inference rates compared to probabilistic mapping algorithms and demonstrate reliability compared to the most similar learning-based comparison, Semantic MapNet [16]. We also present real-world results in an open-source ROS node on challenging off-road data.

This paper is built upon our previous work [17], where we established a new data set and neural network for semantic mapping in dynamic scenes called MotionSC. However, MotionSC was unable to bridge the sim-to-real gap from the CarlaSC simulated data set to real-world Semantic KITTI [2] and required accurate 3D ground truth maps, which are often unavailable in the real world.

Compared to our previous conference paper on Convolutional Bayesian Kernel Inference [19], we accelerate the mapping algorithm by modifying the voxel storage to a sliding local window which discards voxels outside of the local boundaries. Additionally, we add studies on the reliability of ConvBKI which compare ConvBKI with end-to-end mapping network Semantic MapNet on the sim-to-real transfer between Carla SC and Semantic KITTI. We also add more in-depth ablation studies on the effect of dynamic objects, noisy input, and the run-time of local mapping compared to global map storage. Lastly, we gather a real-world perceptually challenging off-road dataset and study the ability of ConvBKI to transfer from off-road driving dataset RELIS-3D [20] to our new dataset, which is unrepresented by the training set.

In summary, our contributions are:

- i. Novel neural network layer for closed-form Bayesian inference, which combines the best of probabilistic and differentiable programming.
- ii. Real-time explicit probabilistic incorporation of semantic segmentation predictions into a map with quantifiable uncertainty.
- iii. Comparison of the proposed ConvBKI approach against the state-of-the-art learning-based and probabilistic semantic mapping methods.

B. Outline

The remaining content of this article is organized as follows. Section II provides a literature review on the evolution of mapping algorithms and off-road perception. Section III provides an overview of Bayesian Kernel Inference necessary to understand ConvBKI. Section IV proposes a new deep learning layer for real-time incorporation of semantic segmentation predictions within a map, trained only on semantic segmentation data. Section V accelerates the framework with local mapping. Section VI compares ConvBKI against learning and probabilistic approaches in addition to ablation studies for the network design. Section VII concludes with a ROS package for open-source use and evaluates ConvBKI on challenging real-world off-road driving data. The ROS package can be found at https://github.com/UMich-CURLY/BKI_ROS with examples on sequences of data from Semantic KITTI [2] and RELIS-3D [20]. The new off-road data is unable to be released

publicly due to proprietary reasons. Finally, limitations and possibilities for future work are discussed in Section VIII.

II. LITERATURE REVIEW

In this section, we review 3D semantic mapping and the trade-offs between learned and hand-crafted approaches. We also review approaches for off-road perception, and the importance of uncertainty in perceptually challenging environments.

A. Semantic Mapping

Semantic mapping is the task of incorporating semantic labels such as car or road into a geometric world model. Historically, most mapping methods were hand-crafted and mathematically derived. Early semantic mapping algorithms semantically labeled images, then projected to 3D and directly updated voxels through a voting scheme or Bayesian update [4], [21]–[23]. Later semantic mapping algorithms applied further optimization through Conditional Random Fields (CRF), which encourages consistency between adjacent voxels [24]–[26]. Separately, continuous mapping algorithms estimate occupancy through a continuous non-parametric function such as Gaussian processes (GPs) [11], [27]. However, these methods suffer from a high computation load, rendering them impractical for onboard robotics with large amounts of data. For example, GPs have a cubic computational cost with respect to the number of data points and semantic classes [28]. Bayesian Kernel Inference is a method proposed to accelerate continuous mapping by approximating GPs at model selection, which will be discussed in more detail in Section III.

Other works have also explored alternative data-efficient representations such as surfels [6], truncated signed distance functions [7], and meshes [8], [10], although these methods are not exclusive of semantic mapping and can be used as an extension to create high-fidelity semantic maps [10], [29]. While in this work we focus on real-time semantic mapping, our algorithm can be extended to create real-time semantic meshes through methods such as VoxBlox [30].

Many modern approaches to mapping leverage inventions in the deep learning community to learn an efficient, implicit approximation of the world in a lower dimensional latent space. Transformers are large arrays of multi-head attention layers [31] which can be combined to form large language models [32], vision transformers [33], and networks for 3D object detection [34] and mapping [15] among other applications. However, transformers are notoriously difficult to train and suffer un-diagnosable failures when exposed to new data. Other modern approaches to semantic mapping apply recurrent neural networks [16], [35], [36] or spatiotemporal convolution networks [14], [17] to model spatiotemporal dynamics. Spatiotemporal approaches remove the recurrency of other mapping algorithms and instead aggregate past information to operate directly on a finite length temporal dimension of information. While operating on a finite length time dimension is efficient for convolutional networks, the assumption discards valuable memory. Additionally, learning approaches for mapping require large amounts of supervised data which

can be challenging to obtain in real-world dynamic driving environments [2].

Although learning-based approaches have achieved success in minimizing memory and accelerating inference, they still encounter significant challenges. By approximating functions implicitly, there is no notion of when a network will fail, as provided by variance, or the ability to diagnose an error. Additionally, latent space operations are vulnerable to errors when exposed to new data, motivating research into problems such as the sim-to-real gap. On the other side of the spectrum, mathematical hand-crafted approaches provide reliability and trustworthiness at the cost of efficiency. Our approach seeks to find a middle ground by operating on explicit probability distributions in an end-to-end differentiable neural network with modern deep learning layers. Next, we discuss literature in off-road perception which necessitates a balance between efficiency and reliability.

B. Off-Road Perception

Perception in off-road driving is challenging due to less structured environments with more ambiguity and unbalanced data than on-road driving. Roads generally contain a finite set of easily distinguishable classes with a regular geometric pattern, such as the location and shape of roads, vehicles, traffic lights, etc. However, off-road presents more irregular geometric shapes which are not easily segmented and depend on context [20]. Although on-road perception has advanced tremendously due to a plethora of large-scale datasets [2], [37], [38], fewer off-road semantic segmentation datasets are available [20], [39], [40]. RELIS-3D is the largest off-road LiDAR semantic segmentation data set at the time of writing, containing 13,556 labeled point clouds. The labels are unequally distributed, dominated by grass, tree, and bush semantic categories. Uncertainty-aware semantic segmentation [41] significantly outperforms baselines [42] on RELIS-3D [20], emphasizing the importance of uncertainty quantification on challenging unstructured outdoor perception.

Other works have found that geometric information alone is insufficient for off-road autonomy due to ambiguous structure [43]–[45], necessitating more information on traversability through semantic mapping [40]. For example, without semantic context, a bush or tall blade of grass may be confused as impenetrable obstacles. Due to these challenges of off-road navigation, we test the semantic mapping and uncertainty quantification of ConvBKI on a new off-road data set in Section VI. In the next section we introduce background on efficient probabilistic mapping necessary to understand our proposed method.

III. PRELIMINARIES

Semantic Bayesian Kernel Inference (S-BKI) [5] is a 3D continuous semantic mapping framework which builds on the work of Vega-Brown et al. [13] and Doherty et al. [12]. Semantic BKI recursively incorporates semantically labeled points into a voxelized map by considering the relative position between each input point and voxel through a kernel function. Bayesian Kernel Inference (BKI) is an efficient approximation

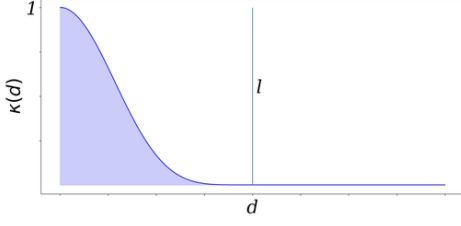


Fig. 2: Sparse Kernel Function. $\kappa(d)$ has a maximum value of 1 at $d = 0$, and decays to 0 by $d = l$. When applied to semantic mapping, points proximal to the voxel centroid have more influence over the semantic label of the voxel than points further from the centroid.

of GPs which leverages local kernel estimation to reduce computation complexity from $\mathcal{O}(N^3)$ operations to $\mathcal{O}(\log N)$, where N is the number points. In contexts such as mapping, there may be hundreds of thousands of points, rendering GPs impractical.

For supervised learning problems, our goal is to identify the relationship $p(y_*|x_*, \mathcal{D})$ from a sequence of N independent observations $\mathcal{D} = \{(x_1, y_1), \dots, (x_N, y_N)\}$ for query point x_* with label y_* . In semantic mapping, the input data is a series of semantically labeled points with positions x_i and semantic observations y_i , and $*$ represents the query voxel. Observations are associated with latent model parameters θ_i which define the likelihood $p(y_i|\theta_i)$. In the case of semantic mapping, we wish to infer the categorical parameters θ_i , such that the probability of observing a one-hot encoded semantic category y_i may be written as:

$$p(y_i|\theta_i) = \prod_{c=1}^C (\theta_i^c)^{y_i^c}. \quad (1)$$

As we observe 3D semantic point-pairs (x_i, y_i) , our goal is to learn the latent target parameters θ_* which define the Categorical posterior of each voxel $*$ with centroid x_* . Note that the voxel posterior $p(\theta_*|x_*, \mathcal{D})$ is dependent on observations continuously scattered around the voxel, requiring an update which considers the relative position of input data.

Vega-Brown et al. [13] propose a solution to relate the extended likelihood $g(y_i) = p(y_i|\theta_*, x_i, x_*)$ to the likelihood distribution $f(y_i) = p(y_i|\theta_*)$ through a positional kernel $k(x_i, x_*)$. This relationship generalizes local kernel estimation to Bayesian inference for supervised learning. In particular, they relate the extended likelihood and likelihood by showing that the maximum entropy distribution g satisfying $D_{KL}(g||f) \geq \rho(x_*, x_i)$ has the form $g(y_i) \propto f(y_i)^{k(x_*, x_i)}$, where k is a kernel function and D_{KL} is the Kullback-Leibler Divergence. In this case, $\rho : X \times X \rightarrow \mathbb{R}^+$ is some function which bounds information divergence between the likelihood distribution $f(y_i) = p(y_i|\theta_*)$ and the extended likelihood distribution $g(y_i) = p(y_i|\theta_*, x_i, x_*)$. Concretely,

$$p(y_i|x_i, \theta_*, x_*) \propto p(y_i|\theta_*)^{k(x_i, x_*)}. \quad (2)$$

Functions k and ρ have an equivalence relationship, where each is uniquely determined by the other. The only requirements are that

$$k(x, x) = 1 \forall x \quad \text{and} \quad k(x, x_*) \in [0, 1] \forall x, x_*. \quad (3)$$

As described above, in semantic mapping θ_* describes a Categorical distribution over the semantic categories C . Therefore, the extended likelihood spatially relates semantic observations (x_i, y_i) to voxel latent parameters θ_* using a kernel function k . This formulation is especially useful for likelihoods $p(y_i|\theta_*)$ chosen from the exponential family, as the likelihood raised to the power of $k(x_*, x)$ is still within the exponential family.

Doherty et al. [12] then apply the BKI kernel model to the task of occupancy mapping. In occupancy mapping, occupied points are measured by a 3D sensor such as LiDAR, and free space samples can be approximated through ray tracing. Measurement $x_i \in \mathbb{R}^3$ represents a 3D position with corresponding observation $y_i^c \in \{0, 1\}$ where $(y_i^1 = 1)$ for free space measurements and $(y_i^2 = 1)$ for occupied measurements at the end of a ray. The voxel occupancy can be modeled by a Bernoulli distribution over θ_* . Adopting the conjugate prior distribution $\text{Beta}(\alpha_*^1, \alpha_*^2)$ over voxel parameters θ_* yields a closed-form update equation by applying Bayes' rule:

$$p(\theta_*|x_*, \mathcal{D}) \propto p(\mathcal{D}|\theta_*, x_*)p(\theta_*|x_*), \quad (4)$$

which can be written in terms of the extended likelihood as

$$p(\theta_*|x_*, \mathcal{D}) \propto \left[\prod_{i=1}^N p(y_i|x_*, x_i, \theta_*) \right] p(\theta_*|x_*). \quad (5)$$

Applying the relationship between extended likelihood and likelihood defined in Eq. (2), the posterior can be simplified to

$$p(\theta_*|x_*, \mathcal{D}) \propto \left[\prod_{i=1}^N p(y_i|\theta_*)^{k(x_i, x_*)} \right] p(\theta_*|x_*). \quad (6)$$

Due to the conjugate relationship between the prior $p(\theta_*|x_*)$ and likelihood $p(y_i|\theta_*)$ defined in Eq. (1), the update yields a closed-form solution where:

$$p(\theta_*|x_*, \mathcal{D}) \propto \left[\prod_{i=1}^N \left[\prod_{c=1}^C (\theta_*^c)^{y_i^c} \right]^{k(x_i, x_*)} \right] \prod_{c=1}^C (\theta_*^c)^{\alpha_*^c - 1}, \quad (7)$$

which simplifies to

$$p(\theta_*|x_*, \mathcal{D}) \propto \prod_{c=1}^C (\theta_*^c)^{\alpha_*^c + \sum_{i=1}^N y_i^c k(x_i, x_*) - 1}. \quad (8)$$

Since the posterior is proportional to the conjugate prior distribution $\text{Beta}(\alpha_*)$, the update at time-step t for voxel $*$ can be written in closed-form as

$$\alpha_{*,t}^c = \alpha_{*,t-1}^c + \sum_{i=1}^{N_t} k(x_*, x_i) y_i^c. \quad (9)$$

The equation provides a closed-form method for updating the belief that voxel $*$ is occupied or free given observed measurements and samples of free space. The un-normalized posterior $\text{Beta}(\alpha_{*,t}^1, \alpha_{*,t}^2)$ is a weighted summation of nearby inputs points, where influence of nearby points is calculated

through the kernel function. From the un-normalized voxel concentration parameters α_* , the normalized expectation and variance of voxel parameters θ_* are defined as

$$\eta_* = \sum_{j=1}^C \alpha_*^j, \quad \mathbb{E}[\theta_*^c] = \frac{\alpha_*^c}{\eta_*}, \quad \mathbb{V}[\theta_*^c] = \frac{\frac{\alpha_*^c}{\eta_*} (1 - \frac{\alpha_*^c}{\eta_*})}{1 + \eta_*}. \quad (10)$$

Doherty et al. [12] propose to use a sparse kernel function [46] which simplifies kernel calculation to a function of the Euclidean distance between two points. For two points x_i and x_* , the sparse kernel is calculated as

$$\kappa(d) = \begin{cases} \sigma_0 [\frac{1}{3} (2 + \cos(2\pi \frac{d}{l}) (1 - \frac{d}{l}) + \frac{1}{2\pi} \sin(2\pi \frac{d}{l}))], & \text{if } d < l \\ 0, & \text{else} \end{cases} \quad (11)$$

where $d := \|x_i - x_*\|_2$ and l is the kernel length. The variables σ_0 is used to scale the magnitude of the kernel, however $\sigma_0 = 1$ to satisfy Eq. 3. Note that whereas the kernel function $k : \mathbb{R}^3 \times \mathbb{R}^3 \rightarrow [0, 1]$ is defined between two points, the sparse kernel function $\kappa : \mathbb{R}^3 \rightarrow [0, 1]$ is a function of the Euclidean distance between the two points. The sparse kernel is shown in Figure 2, and effectively provides more weight to close points.

Gan et al. [5] show that the same approach can be applied to semantic segmentation predictions from neural networks by adopting a Categorical likelihood over θ_* with conjugate prior $\text{Dir}(C, \alpha_0)$. This model is also calculated using Eq. (9), however the number of classes has increased from two to any positive integer C .

Although Semantic BKI has achieved accurate and efficient performance in 3D mapping compared to previous methods, it is still limited in key ways. Firstly, the kernel is hand-crafted, and kernel parameters must be manually tuned. As a result, a single spherical kernel is shared between all semantic classes. Second, the update operation is computationally expensive, as the kernel inference must locate all nearby voxels and does not leverage the parallelization of modern GPU architectures. In Section IV we show how ConvBKI improves upon these shortcomings.

IV. CONVOLUTIONAL BAYESIAN KERNEL INFERENCE

In this section, we propose a novel neural network layer, Convolutional Bayesian Kernel Inference (ConvBKI), intended to accelerate and optimize S-BKI. Compared to S-BKI, ConvBKI *learns* a unique kernel for each semantic class and generalizes to 3D ellipsoids instead of restricting distributions to spheres. We introduce the layer and incorporate it into an end-to-end deep neural network for updating semantic maps. A figure summarizing our approach is shown in Fig. 3.

A. Layer Definition

We build a faster, trainable version of Semantic BKI based on the key observation that Eq. (9) may be rewritten as a depthwise convolution [47]. We find that the kernel parameters

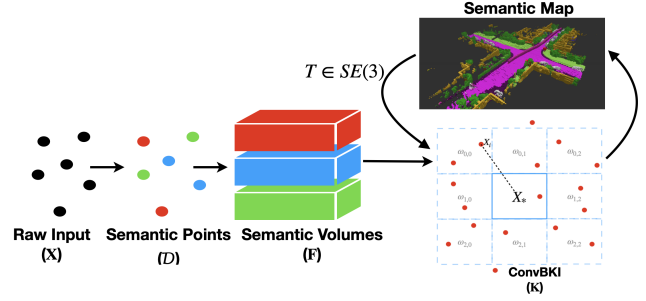


Fig. 3: Structural overview of ConvBKI. Input 3D points are labeled with semantic segmentation predictions by a semantic segmentation neural network, and are voxelized into a dense voxel grid (\mathbf{F}) where each voxel contains the summation over the semantic segmentation predictions of all points within the voxel. A 3D depthwise separable convolution performs the Bayesian update in real-time, where the posterior semantic voxel grid is equal to the convolution of \mathbf{F} plus the prior voxels from the last time-step.

are differentiable with respect to a loss function and, therefore, may be optimized with gradient descent. Learning the kernel parameters enables more expressive geometric-semantic distributions and improved semantic mapping performance.

First, we note that semantic segmentation neural networks predict soft-max encoded predictions, and one-hot encoding the measurement y loses information. Intuitively, the closed-form update is a weighted semantic counting model which should be able to consider soft-max observations from segmentation neural networks. Therefore, we propose to consider the likelihood of soft-max prediction y_i as

$$p(y_i | \theta_i) \propto \prod_{c=1}^C (\theta_i^c)^{y_i^c}, \quad (12)$$

where $\sum_{c=1}^C y_i^c = 1$ and $y_i^c \geq 0$. Substitution into Eq. (7) yields the same closed-form update as Eq. (9), and enables consideration of one-hot encoded or soft-max encoded neural network observations. The distribution defined in Eq. (12) is a specific case of the Continuous Categorical distribution proposed in [48], and has also been noted as useful for the BKI framework in [49].

The update operation in Eq. (9) performs a weighted sum of semantic probabilities over the local neighborhood of voxel centroid x_* . This operation can be directly interpreted in continuous space with radius neighborhood operations such as in PointNet++ [50], DGCNN [51], or KPConv [41]. However, in practice these operations are much too slow to compute for hundreds of thousands of camera or LiDAR points due to an expensive k-Nearest Neighbor operation. Instead, we perform a discretized update, where the geometric position of each local point is approximated by voxelized coordinates. Approximation through downsampling is already performed in Semantic BKI [5], and is a common step in real-time mapping literature [10], [30].

Our method updates a dense voxel grid of a fixed size at each time step, with voxel concentration parameters initialized to prior at the first time step. For the prior, we use a small positive value such $\alpha_*^c = \epsilon$ for all voxels and semantic categories. We apply a discretized update to the dense voxel grid of dimension $\mathbb{R}^{D_C \times D_X \times D_Y \times D_Z}$, where D represents

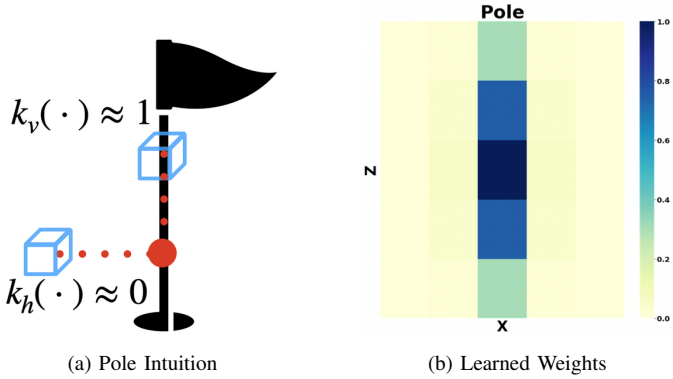


Fig. 4: Illustration of compound kernel motivation. ConvBKI learns a distribution to geometrically associate points with voxels. Whereas a point (red) labeled pole suggests a vertically adjacent voxel (blue) may also be a pole, it does not imply the same for a horizontally adjacent voxel. The image on the right demonstrates a slice of the learned weights for the class pole, which matches the expectations shown on the left.

the dimension of the semantic channel (C) and Euclidean (X, Y, Z) axes.

For each input point cloud, we discretize input points to voxelized coordinates and create a dense semantic volume \mathbf{F} , where input \mathbf{F}_*^c is the sum of the probability mass for semantic category c of all predictions y_i contained in voxel $*$. Let $I(*, i)$ be an indicator function representing whether point x_i lies within voxel $*$. We compute input semantic volume $\mathbf{F} \in \mathbb{R}^{D_C \times D_X \times D_Y \times D_Z}$ as follows.

$$\mathbf{F}_*^c = \sum_{i=1}^N I(*, i) y_i^c \quad (13)$$

For each voxel, the Bayesian update can be calculated as the sum of the prior and a depthwise convolution over input \mathbf{F} , equivalent to a discretized representation of Eq. (9). Let h, i, j be the discretized coordinates of voxel $*$ within \mathbf{F} , and n, o, p be indices within convolution filter $\mathbf{K} \in \mathbb{R}^{D_C \times f \times f \times f}$ where f is the filter size. Then, we can write the update for a single semantic channel of voxel $*$ as

$$\alpha_{*,t}^c = \alpha_{*,t-1}^c + \sum_{n,o,p} \mathbf{K}_{n,o,p}^c \mathbf{F}_{h+n,i+o,j+p}^c \quad (14)$$

where integer indices $n, o, p \in [-\frac{f-1}{2}, \frac{f-1}{2}]$. Note that this is the equation for a zero-padded depthwise convolution, where dense 3D convolution is performed at each voxel in the feature map, with a unique convolution filter \mathbf{K}^c per semantic category c . As a result, this operation may be accelerated by GPUs and optimized through gradient descent since the convolution operation is differentiable.

While it is possible to learn an individual weight for each position and semantic category in filter \mathbf{K} , we found that restricting the number of parameters through a kernel function increases the ability of the network to quickly learn generalizable semantic-geometric distributions. Following [5] and [12], we use a sparse kernel [46] as our kernel function since the sparse kernel fulfills the requirements listed in Eq. (3). Additionally, the sparse kernel is differentiable so that a partial derivative of the loss function with respect to the kernel parameters can be calculated. A plot of the sparse kernel

function is included in Fig. 2 for reference and is described mathematically in Eq. (11), where the parameters are kernel length l , and signal variance σ_0 . Note that for Eq. (3) to remain valid, σ_0 must be 1, leaving only one tune-able parameter for the kernel function.

Effectively, kernel \mathbf{K} is a weight matrix where each weight represents a semantic and spatial likelihood of points being correlated around the query voxel. Therefore, we propose to learn a unique kernel k^c for each semantic category and assign convolutional filter weights to \mathbf{K} by discretizing the kernel function. For a filter of dimension f and resolution Δr , the convolution parameters $\mathbf{K}_{n,o,p}^c$ at filter indices n, o, p can be calculated by evaluating kernel function k^c at the offset of position n, o, p relative to the filter center as follows.

$$\mathbf{K}_{n,o,p}^c = k^c \left(\mathbf{0}, \Delta r \left(\frac{f-1}{2} - \begin{bmatrix} n \\ o \\ p \end{bmatrix} \right) \right) \quad (15)$$

To accommodate complex geometric structures of real objects, we also propose a compound kernel [52, Ch. 4] computed as the product of a sparse kernel over the horizontal plane (κ_h) and vertical axis (κ_v). The compound kernel is defined as

$$k^c \left(\begin{bmatrix} x_i^1 \\ x_i^2 \\ x_i^3 \end{bmatrix}, \begin{bmatrix} x_*^1 \\ x_*^2 \\ x_*^3 \end{bmatrix} \right) = \kappa_h^c \left(\left\| \begin{bmatrix} x_i^1 - x_*^1 \\ x_i^2 - x_*^2 \end{bmatrix} \right\|_2 \right) \kappa_v^c \left(\|x_i^3 - x_*^3\|_2 \right), \quad (16)$$

where the integer superscript represents the Euclidean axes.

Intuitively, ConvBKI treats the output of a semantic segmentation neural network as sensor input and learns a geometric probability distribution over each semantic class. Semantic classes have different shapes, where classes such as poles are more vertical, and classes such as road have influence horizontally. This idea is visualized in Fig. 4, where points vertically adjacent to a pole have a higher weight than points horizontally adjacent to a pole.

Beyond semantic labels, some important considerations for applying ConvBKI include free space measurements and dynamic environments. In our experiments, we focus on semantic quality and do not incorporate free space observations. Instead, we filter free space voxels through sampling by only visualizing voxels with $\eta_* > 0.1$, where a voxel with a measurement would have $\eta_* \geq 1$. In applications where explicit free space is desired, explicit free space measurements can be incorporated into the semantic map by sampling along each ray at fixed distance intervals [12], and treating the sampled free space points as another semantic category when performing the update [5]. Free space sampling is incorporated into our software, and inference results of free-space sampling are included in Section VI-C.

Another important detail to note is that ConvBKI does not draw distinction between static and dynamic objects, leaving artifacts in the map as objects move as traces regardless of how free space is calculated [18]. This concept is demonstrated in Section VI-B, where we purposefully evaluate ConvBKI in highly dynamic scenes without any special consideration paid to dynamic objects. Methods exist which propose to either

separately track dynamic objects from the static world, or use velocity information to decay the occupied belief [18].

B. Training

One key advantage of ConvBKI over alternative learning-based mapping approaches is the cost of data. ConvBKI seeks to find the maximum entropy distribution g , which spatially relates query point i with voxel $*$ as $g(y_i) = p(y_i|\theta_*, x_i, x_*)$. Therefore, the network is learning a function to improve semantic segmentation through the distribution g recurrently. As a result, ConvBKI only requires semantic segmentation ground truth as opposed to a 3D map since the learning problem is framed as spatially and temporally integrating semantic predictions to best perform semantic segmentation of point x_i . In contrast, other mapping approaches rely upon a ground truth semantic map, which is difficult to obtain in the real world [16], [17].

To train ConvBKI, we formulate the learning problem as a semantic segmentation task given the past \mathcal{T} frames of semantic point clouds labeled using an off-the-shelf semantic segmentation network. 3D points are labeled by an off-the-shelf semantic segmentation network to form $\mathcal{D}_{t-\mathcal{T}:t}$, then transformed to the current frame T_t .

Next, a map of size $\mathbb{R}^{D_C \times D_X \times D_Y \times D_Z}$ is initialized to the prior distribution. The prior distribution assigns a small positive value ϵ to each Dirichlet parameter α so that each voxel is initialized with a high variance and uniform expected value following Eq. (10). Experimentally, we set $\epsilon = 1e^{-6}$ to enforce a weak prior. Points with semantic labels are then grouped into voxels to form discretized measurement \mathbf{F}_t through Eq. (13). Finally, the ConvBKI filter is applied to calculate the posterior distribution per-voxel. The filter is trained as a maximum-likelihood estimator by maximizing the likelihood of each point in the final input point cloud. Note that the likelihood of each semantic category c per voxel $*$ is calculated using Eq. (10) which normalizes the probability of each semantic category, $\mathbb{E}[\theta_*^c]$, by the sum of the un-normalized concentration parameters, η_*^c . Since ConvBKI directly outputs semantic likelihoods, it is trained using the negative log-likelihood loss instead of cross-entropy loss [53].

V. LOCAL MAPPING

The input to ConvBKI at each time step is a point cloud voxelized into a dense local map L_t with semantic predictions \mathbf{F}_t . As the number of voxels in the semantic map increases with time, the cost of constructing L_t grows exponentially. While ConvBKI may run at speeds of over 100 Hz depending on the size of the map, the bottleneck quickly becomes the query operation from the global map. Therefore, we propose to maintain a dense local map in memory and discard voxels outside of the local region similar to [29], [54]. In practice, a dense world model may be maintained in memory and integrated with a lower definition global model [55], similar to planning through a local and global planner. Additionally, simultaneous localization and mapping (SLAM) algorithms have already achieved high inference rates and accuracy [56], so we prioritize local semantic scene understanding over loop

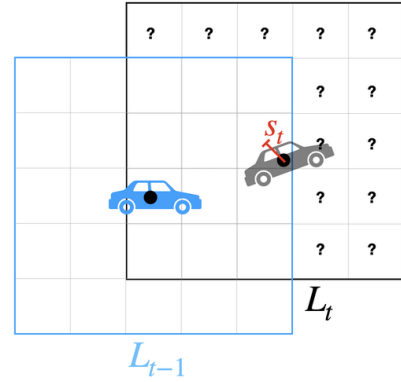


Fig. 5: Toy 2D illustration of local mapping method. When the ego vehicle moves from time $t-1$ to a new location at time t , the ego-vehicle position is approximated to the nearest voxel centroid. The local map L_{t-1} shifts to the new map center and discards voxels outside of L_t . New voxels are initialized to prior, represented by the question mark. Since the center of map is discretized, an offset represented by s_t in red is maintained to transform input point clouds to the map frame. This method prioritizes local semantic scene understanding over loop closure, as voxels will be initialized to prior when re-visited.

closure which may result from global maps. An illustration of our local mapping approach is shown in Fig. 5.

While transforming the local map to match the continuous rotation and translation of the ego-robot is an intuitive approach, it requires approximations of the concentration parameters. We propose to exactly shift the local map by multiples of the voxel resolution, while rotating and translating the input point cloud data to map the local map frame. 3D input may be easily rotated and translated to match the local map coordinate frame without approximations, leading to lossless propagation. Lossless propagation is especially important for explicit probability distributions, as compounding approximations lead to blurring of the map. This approach is applied in real-world robots such as quadrupeds [57] and is shown in Fig. 5.

Our approach consists of several simple steps, optimized for minimal overhead. At each time-step a new pose $T_t \in \text{SE}(3)$ is provided by a localization algorithm. From this pose, we calculate the transformation to the starting pose T_0 :

$$T_{t,0} = T_0^{-1}T_t. \quad (17)$$

Each relative transformation $T_{t,0}$ is composed of a rotation $R_t \in \text{SO}(3)$ and a translation $\psi_t \in \mathbb{R}^3$. The translation between time $t-1$ and t is then calculated in voxel coordinates, dependent on voxel resolution Δr :

$$\nu_{t,t-1} = \left\lfloor \frac{\psi_t}{\Delta r} \right\rfloor - \left\lfloor \frac{\psi_{t-1}}{\Delta r} \right\rfloor, \quad (18)$$

where $\lfloor \cdot \rfloor$ indicates the floating point value is rounded to the nearest integer. Last, we allocate a new local map \hat{L}_t of the same shape as L_{t-1} and copy voxel data within the boundaries of \hat{L}_t from L_{t-1} with an index offset of $\nu_{t,t-1}$.

Since the local map was translated by an integer multiple of the voxel resolution, there may be an offset and rotation from ego coordinates to map coordinates. Therefore, the transformation from sensor coordinates to map coordinates is composed of rotation R_t and an offset

$$s_t = \psi_t - \Delta r \left\lfloor \frac{\psi_t}{\Delta r} \right\rfloor. \quad (19)$$

Intuitively, the translation in voxels of the local map is identified by rounding the translation from the initial coordinate frame to the current ego-centric coordinate frame. Since the local map translation is rounded to an integer, an offset s_t must be maintained to transform input point clouds from the ego-centric frame to the local map frame. Points must also be rotated by R_t to match the orientation of the initial position. This approach requires minimal computational overhead and removes any approximation to ensure reliability.

VI. RESULTS

In this section, we evaluate ConvBKI against probabilistic and learning-based approaches on several data sets with both camera and LiDAR sensors. We also present ablation studies demonstrating the effects of design choices and visualizing the intuitive kernels learned by ConvBKI.

For all results, we train ConvBKI with a learning rate of 0.007 using the Adam optimizer [58] for one epoch with the weighted negative log-likelihood loss. We initialize the kernel length parameter for each kernel to $l = 0.5$ m and train with the last $\tau = 10$ frames as input. The kernel size is set to $f = 5$ with a resolution of 0.2 m.

Inference rates are reported for the dense voxel update step of ConvBKI on an NVIDIA RTX 3090 GPU. Latencies are measured for a batch size of 1 and a single frame, averaged over 1000 repetitions. Note that the update inference rate is directly correlated with the size of the voxel grid, and the map inference rate also depends on free space sampling and voxel grid construction which are studied in Sec. VI-C. However, even including the latency of the semantic segmentation network and a voxel grid with resolution 0.1 m, ConvBKI maintains real-time inference rates with frequencies greater than 10 Hz.

A. Probabilistic Comparison

First, we compare our network to its most direct probabilistic comparisons on the KITTI Odometry dataset [62] with semantic labels obtained from [59] following the approach of Semantic BKI [5]. The purpose of comparing ConvBKI against probabilistic baselines is to highlight the quantitative advantages of kernel optimization, as well as the accelerated inference rate.

The test set contains a sequence of one hundred consecutive raw RGB images with depth estimated by ELAS [63] and

TABLE I: Results on KITTI Odometry sequence 15 [59]. As the degree of kernel expressivity increases, mIoU of ConvBKI also improves from 76.5% without optimization to 77.7% with fully optimized per-class compound kernels.

Method	Building	Road	Veg.	Sidewalk	Car	Sign	Fence	Pole	mIoU (%)	Freq. (Hz)
Segmentation [60]	92.1	93.9	90.7	81.9	94.6	19.8	78.9	49.3	75.1	
Yang et al. [61]	95.6	90.4	92.8	70.0	94.4	0.1	84.5	49.5	72.2	
BGKOctoMap-CRF [12]	94.7	93.8	90.2	81.1	92.9	0.0	78.0	49.7	72.5	
S-CSM [5]	94.4	95.4	90.7	84.5	95.0	22.2	79.3	51.6	76.6	
S-BKI (fine)	94.6	95.4	90.4	84.2	95.1	27.1	79.3	51.3	77.2	0.6
S-BKI (0.2m)	92.6	94.7	90.9	84.5	95.1	21.9	80.0	52.0	76.5	
ConvBKI Single	92.7	94.8	90.9	84.7	95.1	22.1	80.2	52.1	76.6	
ConvBKI Per Class	94.0	95.5	91.0	87.0	95.1	22.8	81.8	52.9	77.5	
ConvBKI Compound	94.0	95.6	91.0	87.2	95.1	22.8	81.9	54.3	77.7	44.3

pose estimated by ORB-SLAM [64]. Semantic segmentation predictions for the image were generated by the deep network dilated CNN [60]. Ground truth semantic segmentation labels are available for twenty-five of the one hundred input images test set images from [59]. Each baseline recurrently generates a global map from the stream of time-synchronized pose, depth, and semantic segmentation predictions. The goal of each algorithm is to estimate semantic segmentation labels at the current time step, given all past information. Therefore, baselines are compared by projecting pixels with depth to 3D world coordinates and assigning semantic segmentation to each pixel corresponding to the voxel the 3D projection resides in. Semantic predictions are aggregated over all time-steps in a global semantic map, and mapping accuracy is evaluated by the Intersection over Union (IoU) metric.

Baselines include a CRF-based semantic mapping system [61], BGKOctoMap-CRF [5], [12], S-BKI [5] and Semantic-Counting Sensor Model (S-CSM) [5]. Our most direct comparison is S-BKI, as the goal of our network is to accelerate and optimize the algorithm in a deep learning framework. However, S-BKI uses a resolution of 0.1 meters with point downsampling instead of discretization. In contrast, our method discretizes input point clouds into voxels to simplify the kernel calculation as a convolution. Therefore, to directly evaluate the benefit of optimization, we add a version of S-BKI with a resolution of 0.2 meters and discretization, which we label S-BKI (0.2m). Grid bounds for ConvBKI are set to a range of [-40, -40, -5.0] to [40, 40, 5.0] m. Results are summarized in Table I.

Kernel Optimization: The results of Table I demonstrate the utility of optimizing kernels, as subsequently more expressive kernels achieve higher mIoU. Starting from S-BKI (0.2m) which represents an unoptimized version of ConvBKI with a single filter shared between all classes, performance is slightly better than the segmentation input. Even without optimization, the BKI framework is capable of recurrently integrating information to improve semantic segmentation. With optimization over a single kernel shared by all classes, denoted ConvBKI Single, performance slightly increases from 76.5 to 76.6% mIoU. The small increase is due to only a single parameter being optimized, yielding an un-expressive kernel. When a kernel is optimized for each semantic category (ConvBKI Per Class), mIoU increases substantially from 76.6 to 77.5%. Finally, dividing the kernel for each semantic category into the product of a horizontal and vertical kernel (ConvBKI Compound) again increases performance, this time most notably for the pole class. Intuitively, points vertically adjacent to points on a pole are also likely to be a pole, whereas horizontally adjacent points are not. This idea is supported in Section VI-C, which visualizes the kernels learned by ConvBKI Compound.

Latency: Another notable comparison in Table I is between S-BKI (fine) and ConvBKI Compound. S-BKI (fine) uses a finer resolution of 0.1 m with downsampling instead of discretization of the input point cloud. However, by leveraging differentiability, ConvBKI Compound is able to improve semantic mapping performance to surpass that of S-BKI. Additionally, ConvBKI achieves significant acceleration, with

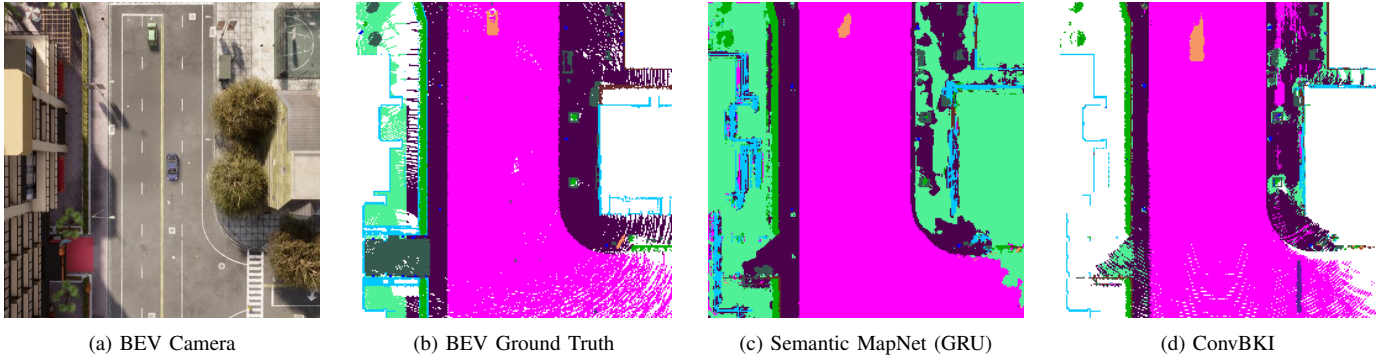


Fig. 6: Comparison of ConvBKI and Semantic MapNet on BEV mapping in CarlaSC. Both networks correctly segment the driveable surface and buildings. Compared with Semantic MapNet, ConvBKI does not make predictions on unseen regions, creating some sparsity in the map. ConvBKI correctly identifies the pedestrian in the bottom right, however leaves a trace as the pedestrian walks along the crosswalk. In contrast, Semantic MapNet omits the pedestrian to achieve a higher quantitative performance.

an inference frequency of 44.3 Hz compared to 2 Hz for SBKI with downsampling and 0.6 Hz without. In summary, ConvBKI achieves an accelerated inference rate and optimized performance compared to probabilistic baselines.

B. Neural Network Comparison

Next, we compare ConvBKI Compound against several types of memory for recurrent neural networks on the task of Bird’s Eye View (BEV) semantic mapping. The purpose of this evaluation is to isolate the trade-offs of ConvBKI opposed to more expressive algorithms with thousands of times more parameters. Specifically, we compare the ground-truth data used, memory, quantitative performance, inference rate, and reliability when transferring to a different dataset.

In contrast to the global mapping comparison for probabilistic mapping, we evaluate on local BEV semantic mapping for a learning comparison. BEV mapping is more common for learning-based algorithms due to the memory cost of 3D operations, and provides a measure of scene understanding by evaluating the semantic map of ConvBKI from a different perspective.

We chose to compare with Semantic MapNet [16] as it was the most direct open-source learning-based comparison we were able to find. Semantic MapNet is a BEV semantic mapping algorithm which uses features from the last layer of a pre-trained semantic segmentation algorithm as input to per-cell recurrent units which perform a map update in feature space. Semantic MapNet groups points with features from the semantic segmentation network into square cells and applies a

BEV filter which discards all but the highest point in each cell. Each cell with an input point is then updated through a learned memory model, including a Gated Recurrent Unit (GRU) [65], Long Short-Term Memory (LSTM) [66] or a linear layer. At inference time, a series of convolution layers are applied to the latent-space memory to produce a BEV semantic segmentation prediction for each 2-dimensional cell.

Semantic MapNet functions similarly to ConvBKI except for several key differences. First, Semantic MapNet operates in two dimensions using only the highest point in each map cell as input, compared to the 3D representation of ConvBKI which incorporates every input point. Additionally, Semantic MapNet applies a recurrent network which updates each cell using latent features of the input points, whereas ConvBKI applies an explicit Bayesian update on the semantic predictions of the semantic segmentation network.

The input to each ConvBKI and Semantic MapNet is from the same semantic segmentation network, Sparse Point-Voxel Convolutional Neural Network (SPVCNN) [67]. SPVCNN is an efficient semantic segmentation for LiDAR point clouds which leverages sparsity to accelerate inference. We follow the training procedure from Semantic MapNet with two modifications, including changing the number of input frames to $\tau = 10$ to match ConvBKI and the batch size to 2 in order to reduce training memory. Semantic MapNet was trained for fifteen epochs over the course of three days. For evaluation, each network is provided the past $\tau = 10$ frames as input and evaluated on the accuracy of the BEV local semantic map produced. Following the evaluation procedure of Semantic MapNet, cells with no input are discarded from evaluation.

1) *Simulated Data*: First, we train ConvBKI and Semantic MapNet on the simulated driving dataset CarlaSC gathered from the CARLA simulator [17]. CarlaSC contains complete 3D scenes with ground truth semantic labels for each voxel at a resolution of 0.2 m and a sensor configuration designed to mimic the popular real-world Semantic KITTI [2] driving dataset. Additionally, CarlaSC is a highly dynamic dataset with moving actors which we use to highlight the artifacts left behind by ConvBKI as discussed earlier in Section IV.

Ground truth BEV semantic maps are obtained by ray-tracing within each ground truth 3D map until an occupied voxel is found, and assigning the BEV semantic label to match

TABLE II: Comparison with learning-based memory on simulated data. Metrics are reported for the isolated recurrent networks, excluding the input segmentation network.

Method	GRU	LSTM	Linear	ConvBKI
Map Dimension	2	2	2	3
Memory (MB)	2765	2917	2707	2497
Parameters	1.96 M	2.04 M	1.74 M	22
mIoU (1 epoch)	24.0	25.3	24.7	26.7
mIoU (15 epoch)	35.4	30.9	32.6	N/A
Freq. (Hz)	111.0	107.4	116.8	176.6
Data	BEV	BEV	BEV	Segmentation

the first occupied voxel. If no occupied voxel is encountered, the 2D cell is ignored during training and evaluation. We compare ConvBKI and Semantic MapNet on the test set of CarlaSC, which includes three scenes of 3 minutes in length sampled at 10 Hz. We use the remapped CarlaSC labels which contain 11 semantic categories: free space, building, barrier, other, pedestrian, pole, road, ground, sidewalk, vegetation and vehicles. Models are evaluated on the mIoU of all semantic predictions in the 5,400 frame test-set. Ground truth scenes are at a resolution of 0.2 meters with local boundaries with respect to the on-board LiDAR sensor of (-25.6, -25.6, -2.0) meters to (25.6, 25.6, 1.0) meters.

Efficiency: Key metrics on the test set of CarlaSC are summarized in Table II. Compared to the learning-based baselines, ConvBKI builds and maintains a *three dimensional* map using *less memory* than the baselines take to operate in two dimensions. This difference is because the ConvBKI update is explicitly defined to achieve optimal efficiency by building off prior probabilistic mapping works. ConvBKI only requires two parameters for each semantic category to learn the geometry, for a total of 22 parameters. In contrast, the learning-based memory models contain nearly a hundred thousand times more parameters due to purely implicit operations. While implicit operations potentially enable more expressive functions to be learned, they are also less reliable as demonstrated in the next section. Additionally, ConvBKI achieves a quicker inference rate despite operating in 3-D due to maximally efficient operations. Note that the inference rate is different from Table I due to the change in grid size.

Training: Another advantage of ConvBKI is the form of training data. ConvBKI trains on semantic segmentation ground-truth data which is more easily obtained than ground-truth BEV maps. Obtaining accurate ground-truth maps in the real world is a difficult challenge due to the presence of dynamic objects [17], requiring multiple viewpoints for accurate data. Semantic segmentation ground-truth data is more easily obtained, as evident by the vast number of on-road driving datasets.

The primary trade-off of ConvBKI is highlighted by the

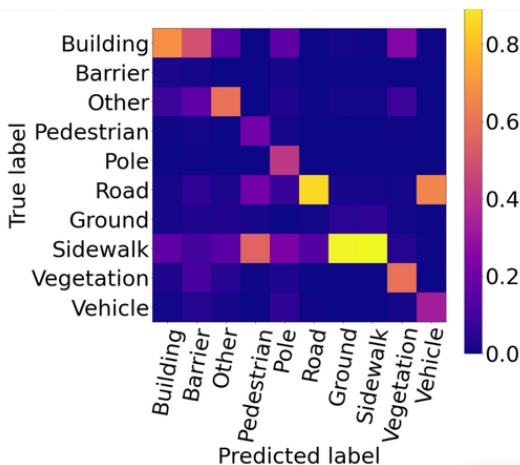


Fig. 7: Confusion matrix of ConvBKI compound on simulated BEV mapping, normalized over the predictions. ConvBKI leaves traces from dynamic objects, causing sidewalk to be incorrectly labeled as pedestrian and road as vehicle.

BEV semantic segmentation performance. ConvBKI trains quickly in less than an epoch, but is gradually outperformed by implicit memory models as they train over the course of several days. With nearly a hundred thousand times more parameters, Semantic MapNet is able to outperform ConvBKI on BEV segmentation. ConvBKI obtains a higher mIoU of 26.7% after one epoch of training compared to the implicit memory models, but is surpassed by the GRU after several days of training which obtains an mIoU of 35.4%.

Dynamic Objects: The lower mIoU of ConvBKI can be explained by examining the precision and recall of each semantic category in Table III. ConvBKI is comparable to Semantic MapNet across most semantic categories, however is penalized on the dynamic classes due to the lack of a dynamic object propagation step considering object motion. ConvBKI recurrently incorporates semantically segmented points into the semantic map without separate treatment of dynamic objects, leaving traces. Measured on BEV mapping, traces cause lower precision of the pedestrian and vehicle classes, as well as lower recall of the sidewalk and road classes. Dynamic methods have already been proposed which fit within the BKI framework [18], however they do not fit into the scope of this paper and we have left the integration of dynamic mapping with our accelerated network as future work.

A confusion matrix of ConvBKI is shown in Fig. 7, and a table of the precision, recall, and IoU per class compared to the fully trained Semantic MapNet is shown in Table III. In general, the performance between the two methods is very similar amongst all classes, with a few exceptions which motivates future work. Some notable columns include pedestrian and vehicle, where ConvBKI has a higher recall yet much lower precision due to artifacts from dynamic objects. In ConvBKI, dynamic objects are treated the same as static objects, which leave traces in the map. When projecting from 3D to 2D, the traces cause ground tiles to be incorrectly labeled as a dynamic object, thereby reducing precision. For the pedestrian class, many pedestrian labels are incorrectly applied to sidewalk tiles, and the same for vehicle predictions to the road class. Note that this is due to the highly dynamic nature of the CarlaSC simulated dataset, and would not be present in static scenes where performance would be higher. The same effect is seen on the road and sidewalk classes, where ConvBKI has a similar precision but lower recall since the BEV projection captures traces left behind by dynamic objects.

Reliability: Another key result is highlighted by the groups of ground classes and building/barrier classes. Our method

TABLE III: Per-class comparison of ConvBKI and SemanticMapNet (SMNet) on simulated data. Key results are bolded. Performance of both algorithms is similar, except on the dynamic object categories. ConvBKI leaves traces which penalizes precision of dynamic objects and recall of ground surfaces.

Method	Metric	Building	Barrier	Other	Pedestrian	Pole	Road	Ground	Sidewalk	Vege.	Vehicle
ConvBKI	Prec.	59.8	1.3	51.5	10.5	29.3	85.1	5.2	91.2	57.4	21.3
	Rec.	67.7	1.3	43.5	31.5	50.2	87.2	6.06	92.6	60.3	56.5
SMNet	Prec.	28.3	18.5	20.4	72.4	44.9	88.7	42.4	41.3	85.4	69.4
	Rec.	34.2	13.2	26.0	35.4	55.3	98.1	22.5	67.7	88.9	53.0
ConvBKI	IoU	23.8	1.2	17.1	10.1	21.5	76.8	4.9	39.7	52.3	19.4
	IoU	29.4	1.2	19.4	20.0	35.7	85.8	5.0	64.2	56.1	37.7

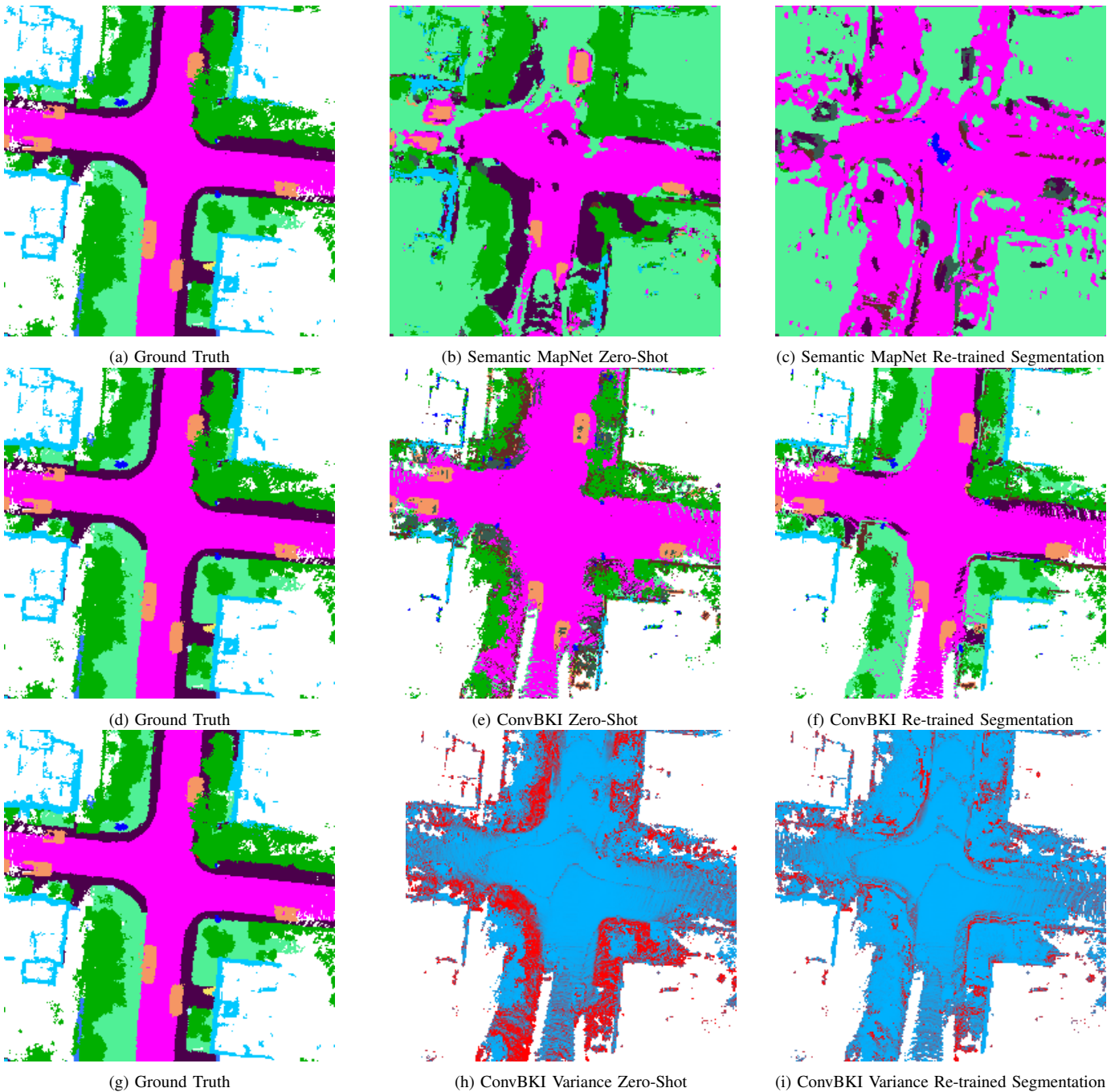


Fig. 8: Comparison of ConvBKI and Semantic MapNet on semantic mapping transferring from simulated to real data. Regions in white are unobserved by ConvBKI or in the ground truth. ConvBKI is generally successful, however suffers issues on ground segmentation due to LiDAR input. Ground segmentation is also a challenge on the simulated data, as the input semantic segmentation network is sometimes unable to distinguish sidewalk from road.

is ultimately reliant on the input semantic predictions of the semantic segmentation network. Although ConvBKI may smooth out noise, it is incapable of resolving errors where the segmentation predicts temporally consistent incorrect labels. However, the predictable nature of our algorithm may also be seen as an advantage for improved reliability, as observed occupied voxels will always be modeled within our map, even if the label is wrong.

In Fig. 6, a pedestrian walks along the bottom right on a sidewalk. ConvBKI leaves a trace as it does not yet handle dynamic objects, while Semantic MapNet completely removes the pedestrian. This occurs because deep neural networks are

constructed to optimize a specific metric, in this case mIoU, and may hallucinate to create the best performance at the cost of wrong predictions. In contrast, ConvBKI is heavily penalized on the pedestrian and road class for leaving traces. In safety-critical environments such as self-driving cars, however, it is vital to have a trustworthy system which does not collide with pedestrians or other difficult semantic categories, even if the semantic label is incorrect.

2) *Transfer to Real Data*: Next, we highlight the reliability of our network by transferring from simulated CarlaSC to real-world Semantic KITTI. The goal of this section is to evaluate the uncertainty quantification ability of ConvBKI under stress,

as well as the ability of ConvBKI and Semantic MapNet to translate to a different dataset. Due to the limited parameters and explicit domain knowledge embedded within ConvBKI, it is more reliable than purely implicit memory models and capable of more successfully bridging the sim-to-real gap. Additionally, the explicit operations of ConvBKI allows the semantic segmentation network to be directly swapped for one trained on the specific environment.

First, we compare ConvBKI with Semantic MapNet on the validation set of Semantic KITTI without any re-training of the semantic segmentation network or memory network. Next, we re-train the semantic segmentation network without modifying the memory networks to demonstrate the advantage of the explicit domain knowledge of ConvBKI.

Ground-truth data is obtained from the Semantic KITTI Semantic Scene Completion task by re-constructing voxels over the same bounds as CarlaSC. Normally, the Semantic Scene Completion task excludes the scene posterior to the ego vehicle, so we re-run the Semantic KITTI voxelizer to ensure BEV map bounds match between CarlaSC and Semantic KITTI voxels. Next, we obtain BEV ground truth from the voxelized scenes following the same procedure as before. ConvBKI and Semantic MapNet are then run directly on Semantic KITTI without any re-training to attempt zero-shot domain transfer from simulation to the real world.

Direct Transfer: ConvBKI is able to more successfully bridge the sim-to-real gap due to explicit domain knowledge which is less susceptible to noise than complex implicit operations. However, error is still introduced by the semantic segmentation network SPVCNN which has not been re-trained. The most common error in ConvBKI is incorrectly labeling other types of ground as road, which is a common systematic error of LiDAR-based semantic segmentation.

Quantitative findings are summarized in Table IV, and predicted BEV maps are shown in Fig. 8. The performance of Semantic MapNet decreases from a mIoU of 35.4% to 23.9% when bridging the sim-to-real gap, whereas ConvBKI has a less significant decrease from 26.7% on simulated data to 24.2% when directly transferring from simulated data to real data. On the new data-set, ConvBKI is able to outperform Semantic MapNet due to a less significant decline in performance. Additionally, when bridging the sim-to-real gap, Semantic MapNet becomes unable to effectively recognize the pedestrian class. As previously stated in Fig. 6, deep learning methods for mapping learn to hallucinate the scene and may unreliably predict the absence of critical objects such as pedestrians to optimize a metric. When transferring

TABLE IV: Quantitative results transferring sim-to-real from CarlaSC [17] to validation sequence of Semantic KITTI [2]. Results for each method are measured once directly from sim-to-real (direct) and once with the semantic segmentation network replaced, indicated by (tran.).

Method	Building	Pedestrian	Pole	Road	Ground	Sidewalk	Vege.	Vehicle	mIoU (%)
ConvBKI (direct)	21.1	17.4	2.8	41.8	19.4	11.4	49.3	30.1	24.2
SMNet (direct)	22.2	3.7	7.0	40.3	22.2	18.4	49.7	27.8	23.9
ConvBKI (tran.)	58.5	23.6	12.3	66.7	64.7	40.5	61.6	45.1	46.6
SMNet (tran.)	0.0	0.0	0.0	18.5	13.7	1.3	3.2	0.3	4.6

from simulation to the real world, pedestrians become significantly more difficult to identify, causing the network to omit pedestrians. In contrast, our method still predicts the presence of pedestrians but is penalized due to leaving artifacts from dynamic objects as previously discussed.

Explicit Knowledge: Another prominent result from Table IV is the difference in performance between ConvBKI and Semantic MapNet when swapping the semantic segmentation network. Since ConvBKI learns explicit distributions over semantic categories, the kernels directly transfer to other semantic segmentation network inputs. The semantic segmentation network may be easily replaced with different weights, in this case, trained on Semantic KITTI. ConvBKI can directly transfer to a new semantic segmentation network since the semantic classes remain the same and may be directly mapped. In contrast, semantic classes in the latent space have no guarantee to correlate.

As expected, the results of ConvBKI significantly improve when the semantic segmentation network performance is improved. Swapping the semantic segmentation network causes the mIoU of ConvBKI to increase from 24.2% to 46.6%. However, since there is no guarantee that features will match in the latent space, the performance of Semantic MapNet decreases when the segmentation network is replaced to 4.6%. These results demonstrate that the world knowledge learned by ConvBKI can readily transfer to other datasets and may be easily improved through the substitution of the semantic segmentation network.

Variance: ConvBI also has the ability to quantify uncertainty, which is highlighted in Fig. 8. Transferring to a new dataset is challenging for the semantic segmentation network, causing high variance in the variance map of ConvBKI. As the segmentation network is swapped, causing an improvement in semantic segmentation, the level of variance in the map decreases. This quality demonstrates the utility of uncertainty quantification, as the variance map provides a measure of confidence for important downstream tasks such as traversability used in path planning. When exposed to new data, the uncertainty can provide warnings to proceed cautiously or stop operation in order to avoid costly errors.

C. Ablation Studies

For our final quantitative results, we analyze design choices on the Semantic KITTI dataset. For this study, we re-train ConvBKI on the validation set of Semantic KITTI for a single epoch and evaluate semantic segmentation on the same set in order to study the upper bound of the expected performance gain. We train and test on a voxel grid with bounds of [-20, -20, -2.6] to [20, 20, 0.6] m along the (X, Y, Z) axes, where points outside of the voxel grid are discarded and not measured in the results. Frequencies are again calculated by the average update latency, with additional results on computation time for input discretization and free space sampling.

Resolution: First, we study the effect of voxel resolution on the inference time and mIoU of ConvBKI. We compare ConvBKI with resolutions 0.1, 0.2 and 0.4 m, and a constant filter size $f = 5$ for all models. Table V indicates that

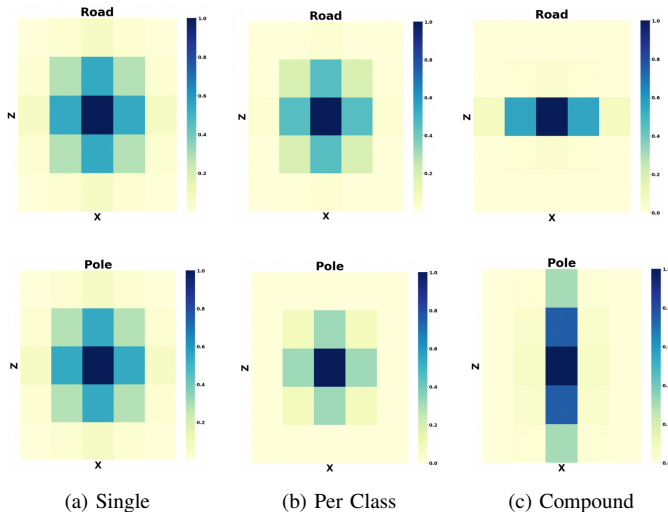


Fig. 9: Illustration of kernels learned by ConvBKI on the road and pole semantic classes, plotted at $\Delta Y = 0$. Adding degrees of freedom increases expressivity by allowing the kernel to learn class-specific geometry.

a finer resolution can increase performance; however, the segmentation difference between 0.2 and 0.1 m resolution is marginal at the cost of greater memory and slower inference. For real-time driving applications, this suggests that 0.2 m resolution may be a strong middle ground. Note that the optimal resolution will vary between applications, such as finer resolution for indoor mobile robots or with LiDAR resolution.

Filter Size: Next, we study the effect of the filter size on inference time and performance summarized in Table VI. While a larger filter size increases the receptive field of the kernel and potentially improves the predictive capability as a result, filter size also cubically increases computation cost. Therefore, identifying a balance between filter size and computational efficiency is important for real-time applications. We study filters of size $f = 3, 5, 7$, and 9 at a resolution of 0.2 m. Table VI demonstrates that filter sizes can improve segmentation accuracy, however, quickly increase run-time. In practice, a filter size of 5 or 7 may be optimal, as a filter size of 9 offers little improvement with a large increase in computational cost. As with resolution, optimal filter size will vary between applications and LiDAR resolution since there will always be a trade-off between performance and latency.

Expressivity: We illustrate the kernels produced by our map to emphasize the underlying operations that match expecta-

TABLE V: Ablation study of voxel resolution on Semantic KITTI sequence 8 for compound ConvBKI with filter size $f = 5$.

Resolution	mIoU (%)	Frequency (Hz)	Mem. (GB)
N/A (Input)	54.6	n/a	n/a
0.4 m	58.2	301.4	2.4
0.2 m	59.3	216.8	2.7
0.1 m	59.0	65.3	5.0

TABLE VI: Ablation study of filter size on Semantic KITTI sequence 8 for compound ConvBKI with resolution 0.2 m.

Filter Size	mIoU (%)	Frequency (Hz)
N/A (Input)	54.6	n/a
$f = 3$	59.0	260.2
$f = 5$	59.3	216.8
$f = 7$	59.5	161.5
$f = 9$	59.6	117.0

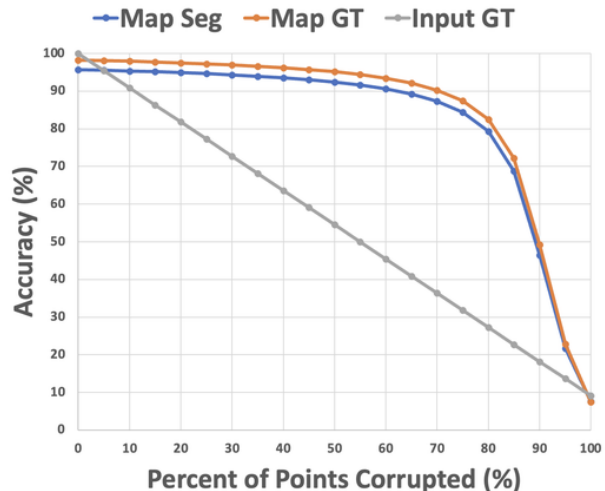


Fig. 10: Ablation study on the smoothing effect of BKI in the presence of noisy segmentation predictions due to aleatoric uncertainty. The horizontal axis indicates the portion of semantic segmentation labels replaced with a random choice from the set of semantic classes. Results are computed as the semantic segmentation accuracy of ConvBKI over every point cloud of Semantic KITTI sequence 8. Two optimal baselines are included for comparison. Input GT refers to the ground truth semantic segmentation without ConvBKI. Map GT refers to the predictions of ConvBKI with ground truth semantic segmentation labels as input. Map Seg refers to the predictions of ConvBKI with neural network predictions as input. Most notably, the segmentation smoothing of ConvBKI minimizes the effect of white noise since spurious predictions are effectively averaged out.

tions of semantic-geometric distributions, such as that poles should be tall and roads should be flat. Fig. 9 demonstrates the kernels learned by variations of the ConvBKI layer for single, per class, and compound kernels. Each variation of ConvBKI improves potential semantic-geometric expressiveness. ConvBKI Single learns one semantic-geometric distribution shared between all classes. However, semantic classes do not share the same geometry in the real world. ConvBKI Per Class adds the capability to learn a unique distribution for each semantic category but is still restricted geometrically. ConvBKI Compound learns a more complex geometric distribution, which can be more expressive for classes with specific shapes.

Robustness: We evaluate the robustness of ConvBKI from aleatoric uncertainty in Fig. 10 compared to optimal baselines. Fig. 10 demonstrates a plot of the percentage of noisy semantic segmentation input labels to ConvBKI versus the overall accuracy. As an optimal baseline, we include a plot of the segmentation accuracy of ground truth labels without mapping (Input GT). As expected, the plot of Input GT begins at an accuracy of 100% and decreases linearly as the percentage of noisy labels increases. Next, we examine the performance of ConvBKI with the ground truth semantic segmentation labels as input (Map GT). Due to compression of points into voxels, the accuracy of Map GT without noise is slightly less than Input GT, but quickly surpasses the performance of Input GT as the amount of noise increases. Next, we plot ConvBKI with predictions from the semantic segmentation neural network as input, denoted Map Seg. ConvBKI performance with semantic segmentation predictions is similar to the optimal baseline with ground truth segmentation, but the gap between accuracies decreases as the amount of noise increases.

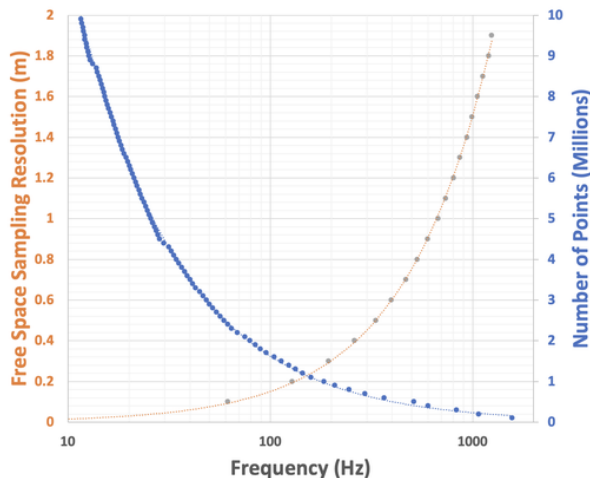


Fig. 11: Frequency of sampling free space points compared to the free space resolution in meters on GPU, and frequency of point insertion into the convolutional input grid compared to the number of points. As free space sampling resolution becomes more fine, the frequency of free space sampling decreases and the number of points grows. Additionally, as the number of points increases, the frequency of input grid construction decreases. While the computation cost of constructing the input grid is dependent on the number of points and free space resolution, the computation complexity of the ConvBKI operation remains constant and is instead dependent on the filter size and map size.

The spatial smoothing operation of ConvBKI reduces the effects of noise at low levels. At high levels of noise, the update step is dominated by noise, causing a sharp decrease in performance. As the amount of noisy predictions increases, the variance of the maximum likelihood prediction also increases due to Eq. (10).

Robustness to epistemic or systemic uncertainty is more difficult for ConvBKI to capture, and is only encapsulated in the variance calculation if objects are more easily distinguishable from different viewpoints. For example, in Fig. 8, road and sidewalk are systematically challenging to distinguish between from LiDAR, but are easier to identify depending on the viewpoint, causing the high variance about misclassified sidewalk regions in ConvBKI zero-shot.

Free Space Sampling: Free space sampling is important for planning frameworks, and is a computationally expensive component of mapping. We sample free space points at fixed intervals along each ray which increases the total number of points. As the number of points increases, the computational complexity required to insert points into the observation grid increases, however the update complexity of the convolutional filter remains constant.

Fig. 11 demonstrates a plot of the frequency of sampling free space points for a point cloud of 120,000 points dependent on the free space sampling resolution between 0 and 2 m. Since the free space sampling equation is vectorized and computed on GPU, the frequency is still well above 10 Hz with a resolution finer than 0.1 m. As the number of points increases due to free space sampling or large data such as fine resolution images, so does the cost of constructing the input for the convolution layer. Fig. 11 also demonstrates the cost of input grid construction compared to the number of input points. Even at 10 million points, grid insertion is still greater

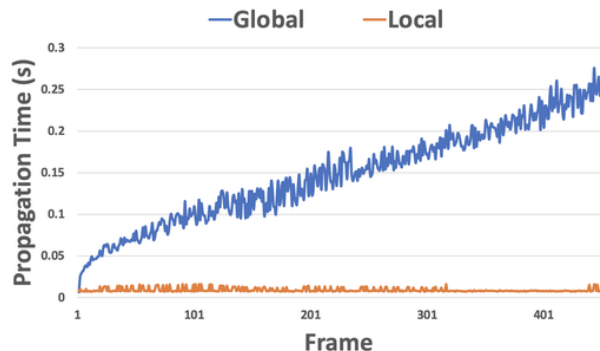


Fig. 12: Propagation time of global and local mapping on Semantic KITTI sequence 8. Since ConvBKI operates on a dense grid, querying the local grid quickly becomes the bottleneck. Local mapping enables accelerated by maintaining the dense local map in memory.

than 10 Hz. Consideration of free space sampling resolution and size of input is necessary for real-time application, as the number of points can rapidly grow depending on parameter choices.

Local Mapping: Last, we study our proposed local mapping method, which we employ for quicker inference. While the update step of ConvBKI may run consistently at high inference rates, the bottleneck of global mapping is querying the dense local region to perform the depthwise separable convolution update. As discussed in Section V, we constrain memory to only maintain the dense local region so the bottleneck may be removed. A comparison of the latency of ConvBKI with global mapping and local mapping is shown in Fig. 12 on the propagation step. Over time, the size of the global map increases, slowing the dense query operation linearly with the number of occupied voxels. In contrast, the local mapping variant has a constant propagation time and removes the bottleneck by restricting map memory to a fixed-size grid which ConvBKI naturally operates on.

VII. OFF-ROAD

Finally, we combine ConvBKI with a semantic segmentation LiDAR network into a Robot Operating System (ROS) package for end-to-end testing. We gather a new off-road data set to study the reliability of our system under perceptually challenging scenarios, including sparse, unstructured data, and when encountering semantic categories previously unseen



Fig. 13: Image from off-road testing. Our dataset includes many semantic categories which are not encountered in RELIS-3D [20], in order to test the reliability of our algorithm in difficult situations.

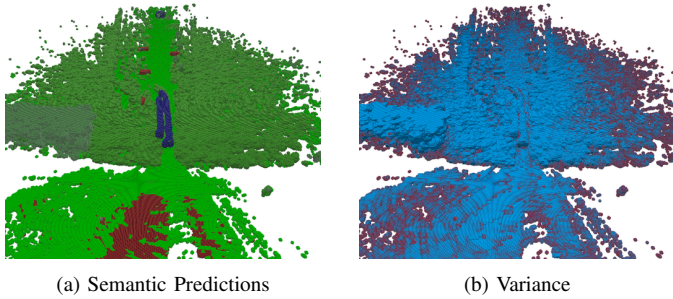


Fig. 14: Expectation and variance of the semantic map on RELLIS-3D dataset. The data set is the largest open-source LiDAR semantic segmentation dataset available at the time of writing but is limited in the number of classes. Typical semantic categories encountered include rubble in red, pedestrian in blue, bush in dark green, and grass in light green, in generally easy scenarios. In contrast, our dataset includes more semantic categories unseen by the robot during training as well as more difficult perceptual conditions.

during training. Although these conditions produce errors, ConvBKI produces quantifiable uncertainty per voxel, which can provide valuable information on when the robot is leaving its operating domain in order to prevent unforeseen errors. An example frame from our data set is shown in Fig. 13, which illustrates the test vehicle as well as the typical scenes encountered in our data set.

The ROS package is available publicly at https://github.com/UMich-CURLY/BKI_ROS. We pre-process poses using an off-the-shelf LiDAR-based SLAM algorithm [56], and perform real-time mapping from recorded ROS bags of real-world scenes. Semantic maps and variance heat maps are published to ROS for visualization as voxel grids. For our semantic segmentation network, we again use SPVCNN [67] for real-time LiDAR semantic segmentation.

A. Training

We train ConvBKI and SPVCNN on the popular off-road driving data set RELLIS-3D [20]. RELLIS-3D contains per-point LiDAR semantic segmentation ground truth labels in a format and size similar to Semantic KITTI [2]. The LiDAR data set contains several imbalanced semantic categories, including grass, tree, and bush, as well as concrete, mud, person, puddle, rubble, barrier, log, fence, and vehicle. The three most common categories are grass, tree, and bush. Scenes in RELLIS-3D are generally well structured, consisting of a robot following a well-defined path. An example semantic and variance map from RELLIS-3D is shown in Fig. 14. On RELLIS-3D, the variance is low except for some regions at the parameters of the map, which are sparsely sampled. Data is captured from a 64-channel LiDAR for high-definition point clouds.

Transferring from RELLIS-3D is particularly challenging as the environments and semantic classes encountered are very different, despite both being off-road data sets. As seen in Fig. 13, there is no clearly defined path in our data set, and a multitude of new classes are encountered. The image displays our test vehicle driving by large trees on an obscure path covered by trees, with buildings to either side. In RELLIS-3D, there is no building class, and categories such as log, fence, or



Fig. 15: Off-road test vehicle equipped with a full sensor suite and ROS.

barrier are rarely seen. These class imbalances and anomalies create unique conditions to test the uncertainty quantification of ConvBKI under stress.

B. Testing

Our data set is gathered with the support of Neya Systems from their full-size military test vehicle, the Polaris MRZR, equipped with a perception sensor suite and Robot Operating System (ROS). The sensor suite includes a 64-channel front-mounted LiDAR, 16-channel rear-mounted LiDAR, stereo camera, IMU, high-precision gyroscope, and GPS. An image of the test vehicle is shown in Fig. 15. Since localization is out of the scope of this paper, we pre-process data using LIO-SAM [56]. LIO-SAM registers point clouds from point cloud and IMU data, providing pose transformations of the LiDAR sensor at each keyframe. After removing points from the front LiDAR within a 1.5 m radius around the LiDAR to avoid points from the ego-vehicle, there are approximately 40 thousand points per point cloud, sampled at 10 Hz. SLAM in off-road driving scenarios is particularly challenging, leaving a small amount of noise in the pose. This noise creates a blurring effect of dense foliage, such as the bushes in Fig. 14.

Data was collected with support from Neya Systems at two of their off-road test facilities with a variety of localization and perception challenges. The first test facility was a large outdoor field with few semantic categories other than bushes and trees. The second, more challenging test facility, was an outdoor paintball facility containing pathways covered by fallen leaves, bare trees, steep hills, buildings, mud, telephone poles, and other difficult features. We recorded several ROS bags over the span of three days of testing. Semantic maps and variance maps were constructed in ROS with the poses from LIO-SAM and filtered front LiDAR.

Images of the raw data as well as semantic and variance maps are shown in Fig. 16 and 17. Fig. 16 includes higher resolution images to demonstrate the sparsity of the data and more vividly portray the maps. In the scene, the vehicle drives along a path with two large buildings on the left, trees and bushes on the right, and another patch of foliage in front.

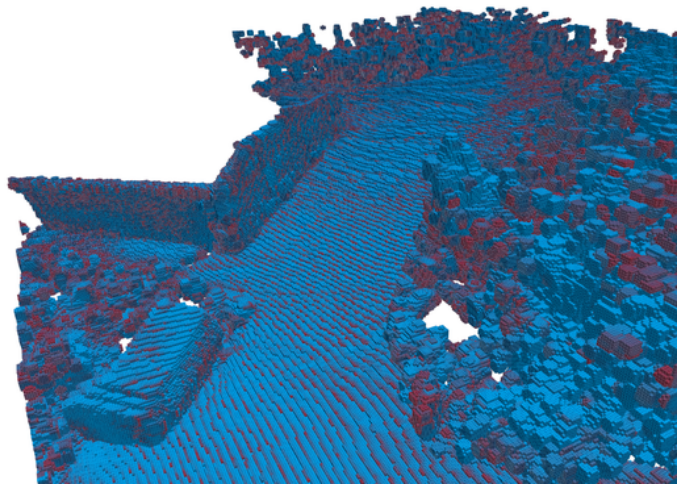
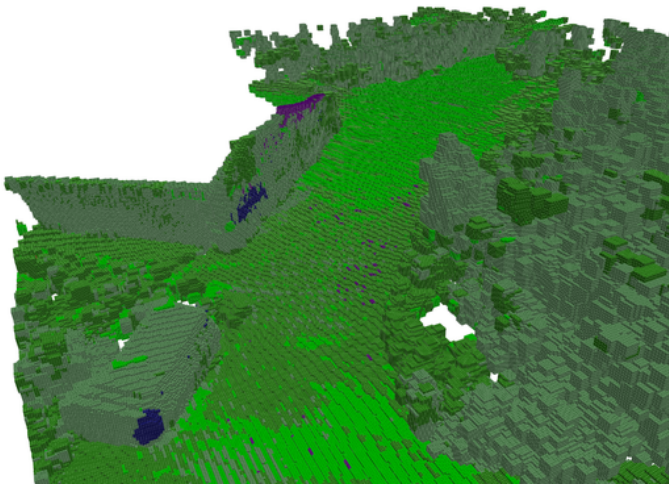
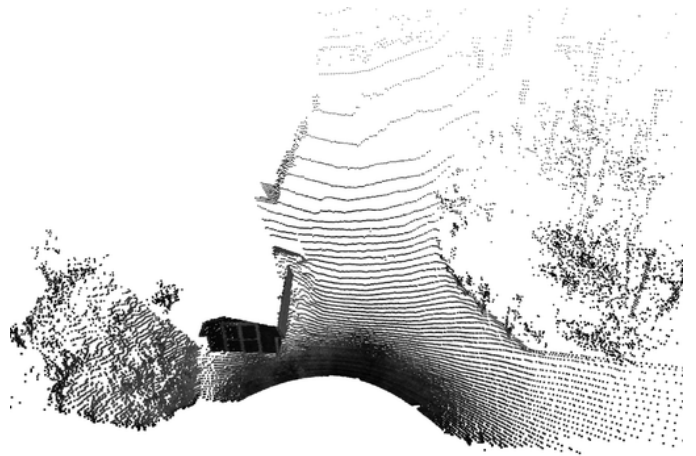


Fig. 16: Example raw data and output from off-road testing. The vehicle drives by two buildings on the left and a collection of trees to the left and front. Trees are in a light brown color, bushes in a dark green, grass in a light green, and pedestrian or pole in dark blue. The small patches of purple in the path are classified as rubble. Since buildings were not encountered during training, points are classified by the closest semantic category the network has seen. Taller portions resemble trees or pedestrians, while flatter objects resemble ground categories such as grass or bush. The point cloud is sparse around clustered bushes and trees, which causes higher variance in dense foliage. Other areas with high variance include misclassified bush and grass on the larger building, as well as the foliage at the end of the path, which has sparsely been observed by the robot.

Despite never encountering the building class, the buildings are labeled with reasonable guesses of mostly tree (brown), with patches of grass (green) or bush (dark green) on the horizontal portions and pedestrian (blue) on the narrow vertical portions. Buildings most closely resemble trees, which are also large un-traversable obstacles. The regions of the large building classified as pole, grass, or bush have high variance and are likely labeled as such due to the sparsity of the point cloud demonstrated in the top right of Fig. 16. At this point, the robot has observed few points from the building and is more prone to making misclassification errors. Other areas with high variance include dense foliage, again due to the difficulty of labeling sparse point cloud data. In particular, voxels at the boundaries between different categories, such as bush and tree or bush and grass, have high variance. The driveable surface is generally segmented correctly with low variance as grass with some bushes and rubble.

Several more examples from our dataset are shown in Fig. 17. In order to include more frames, the images are at a

lower resolution. Each row demonstrates a different example scene from our data, with the raw camera image, point cloud, semantic map, and variance map from left to right. The first row shows an open plain with telephone poles, trees to the sides, and a large building in front. The telephone pole is correctly classified as a pole with low variance, and the building in front is again classified as a tree, with the roof classified as bush. The driveable surface is segmented as grass with some patches of rubble (red) where there are rocks in the path. The patches of rubble have high variance due to sparse point clouds and few rocks, as the network mixes predictions between rubble and grass. Other regions of high variance include the boundaries of dense foliage between bushes and trees.

The next row demonstrates a dense forest scene with tall trees, bushes, and no easily segmented path. The tall trees are sparsely represented in the point cloud, creating misclassifications of trees as poles or bushes. Sparsely covered trees have high uncertainty in the variance map. In contrast, the

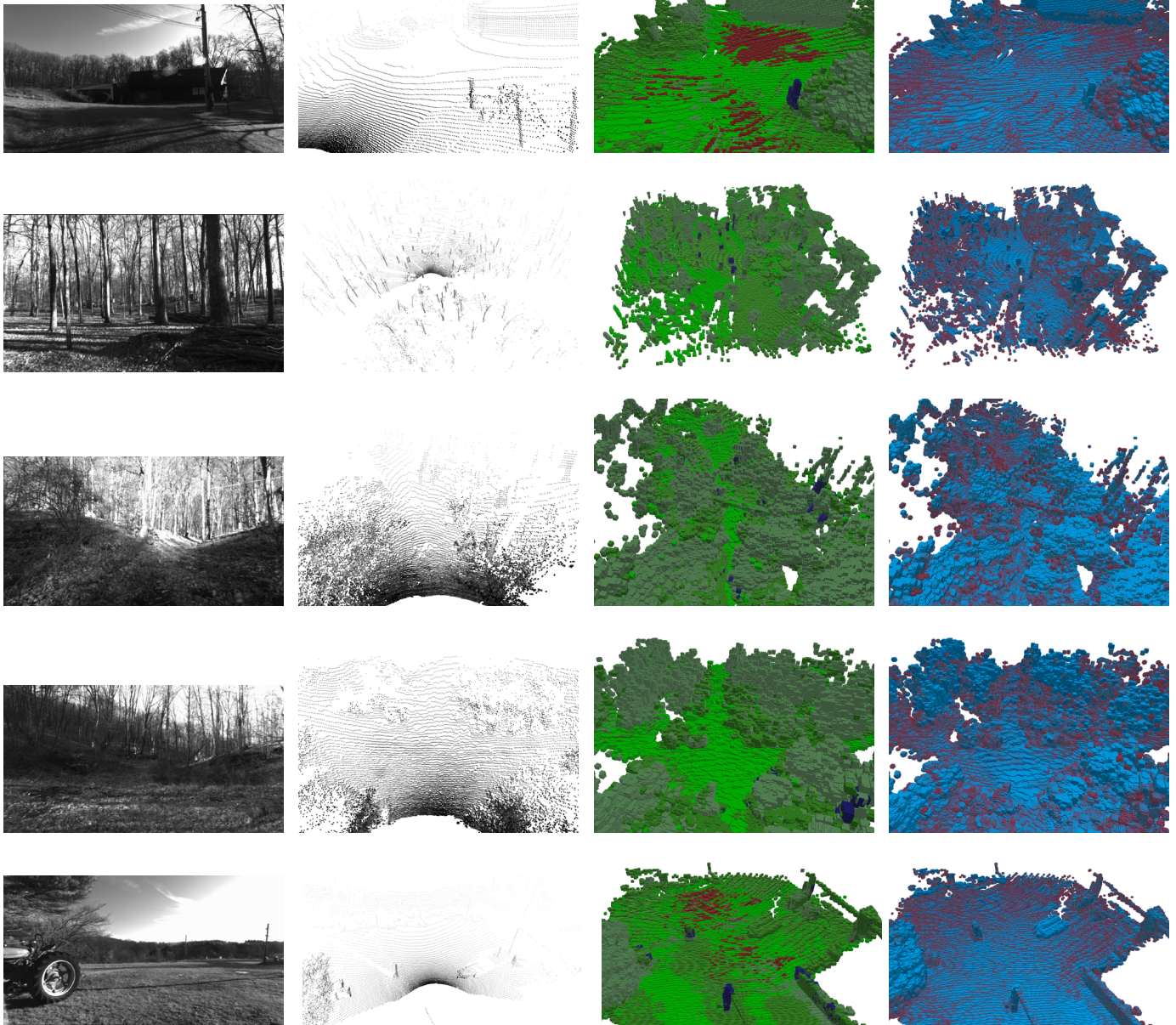


Fig. 17: Raw data and semantic maps from off-road testing. Each row contains a different scene from testing. From top to bottom, the scenes show a building with telephone poles, dense forest, a narrow path through forest, another narrow path through dense forest, and lastly a challenging scene with a tractor, building, and telephone poles. Each row displays the raw sensor data from the camera and point cloud for reference, as well as the semantic and variance maps.

driveable surface of grass and some bushes is easily segmented and has low uncertainty.

The next two rows of images show the vehicle traversing a narrow path through hills surrounded by trees and bushes. Despite the difficult terrain, the driveable surface is correctly segmented with low uncertainty. The dense foliage has high variance at the boundaries of the path and bushes, as well as at mis-classifications of trees as poles.

The final row contains a difficult scene with telephone poles and wires, a tractor, buildings, and more challenging terrain. To the left of the map is a forested region with a tractor to the right side, which is labeled in the semantic map as a cluster of pole (blue), bush, grass and tree, since the semantic

segmentation network is unable to identify the tractor. The raw point cloud of the tractor displays how the wheels may be interpreted as poles or people, and the body of the tractor as bush. In the variance map, the tractor has high uncertainty. The ground is more difficult for the semantic segmentation network to label due to rocks and gravel, resulting in mixed rubble and grass labels in the semantic map and high variance in the variance map. Also included in the map to the right is a small trailer from a truck which is labeled as a bush, and telephone poles with telephone wires. The poles are labeled as pedestrians or trees, while the wires are classified as bushes. The telephone poles and wires have high uncertainty since

they are not captured in the training set.

VIII. CONCLUSION

In this paper, we introduced a differentiable 3D semantic mapping algorithm that combines the reliability and trustworthiness of classical probabilistic mapping algorithms with the efficiency and differentiability of modern neural networks. We quantitatively compared efficiency and accuracy with probabilistic approaches, as well as reliability against a purely implicit deep learning method. To further study the resilience of our model on perceptually degraded driving scenarios, we gathered a new off-road dataset with semantic categories the network had not been trained on and studied the semantic and variance maps.

Overall, ConvBKI is capable of successfully generating maps with meaningful labels despite perceptual challenges, including sparse point clouds, semantic categories outside the training set, and transferring between data sets. By operating explicitly within a differentiable probabilistic framework, ConvBKI maintains reliability on challenging data with quantifiable uncertainty. Additionally, ConvBKI leverages the existing decomposable robotics pipeline, where features are produced from sensor data and recurrently integrated into a map. In contrast, purely implicit deep learning approaches which hallucinate entire scenes can fail unpredictably when exposed to data different from the training set due to a lack of structure and a large number of trainable parameters.

Limitations of ConvBKI: While ConvBKI achieves quick inference rates and expressive, automatically tuned kernels, it has several noteworthy limitations which encourage future work. First, ConvBKI is computationally efficient and parallelizable due to the depthwise convolution operation which applies the same discretized kernel to each semantic category. While the convolution operation accelerates performance, it does not model rotation of points which may be useful in sloped or less structured environments. Second, kernel optimization performs best when the provided training data is similar to the testing data. Training on over-fitted segmentation predictions results in decreased spatial smoothing since the predictions are already near perfect, while training on noisy segmentation predictions results in more expressive spatial kernels. Third, as previously discussed, ConvBKI does not consider the motion of dynamic objects and leaves artifacts in the map.

Future Work: Several avenues for future work exist, and we hope our open source software will encourage other roboticists to extend our method. First, our method can be applied to mesh-based semantic mapping algorithms such as Hydra [29] by replacing the semantic update with ConvBKI. ConvBKI can accelerate the framework and provide voxel-wise variance calculations which may be used to estimate uncertainty at the object level. Another extension is to remove the traces of dynamic objects through a well-formulated network layer which operates explicitly on the Dirichlet distribution concentration parameters such as in Dynamic BKI [18]. Third, while we tested ConvBKI on both camera and LiDAR in this paper, optimal fusion of the two sensors in the BKI

framework remains an open question. Two possibilities are to separately learn kernels for each sensor and add the posterior maps, or to fuse the dense semantic detail from cameras with the accurate geometry of LiDAR before map inference. Finally, ConvBKI may benefit from more expressive extended likelihood distributions which consider the rotation of points instead of only the scale of each category.

REFERENCES

- [1] S. Casas, A. Sadat, and R. Urtasun, "MP3: A Unified Model to Map, Perceive, Predict and Plan," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2021, pp. 14 403–14 412.
- [2] J. Behley, M. Garbade, A. Milioto, J. Quenzel, S. Behnke, C. Stachniss, and J. Gall, "SemanticKITTI: A Dataset for Semantic Scene Understanding of LiDAR Sequences," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2019.
- [3] H. Moravec and A. Elfes, "High resolution maps from wide angle sonar," in *Proc. IEEE Int. Conf. Robot. and Automation*, vol. 2, 1985, pp. 116–121.
- [4] J. McCormac, A. Handa, A. Davison, and S. Leutenegger, "SemanticFusion: Dense 3D semantic mapping with convolutional neural networks," in *Proc. IEEE Int. Conf. Robot. and Automation*, 2017, pp. 4628–4635.
- [5] L. Gan, R. Zhang, J. W. Grizzle, R. M. Eustice, and M. Ghaffari, "Bayesian Spatial Kernel Smoothing for Scalable Dense Semantic Mapping," *IEEE Robot. Autom. Letter.*, vol. 5, no. 2, pp. 790–797, 2020.
- [6] X. Chen, A. Milioto, E. Palazzolo, P. Giguère, J. Behley, and C. Stachniss, "SuMa++: Efficient LiDAR-based Semantic SLAM," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots and Syst.*, 2019, pp. 4530–4537.
- [7] L. Schmid, J. Delmerico, J. L. Schönberger, J. Nieto, M. Pollefeys, R. Siegwart, and C. Cadena, "Panoptic Multi-TSDFs: a Flexible Representation for Online Multi-resolution Volumetric Mapping and Long-term Dynamic Scene Consistency," in *Proc. IEEE Int. Conf. Robot. and Automation*, 2022, pp. 8018–8024.
- [8] M. Herb, T. Weiherer, N. Navab, and F. Tombari, "Lightweight Semantic Mesh Mapping for Autonomous Vehicles," in *Proc. IEEE Int. Conf. Robot. and Automation*, 2021, pp. 6732–6738.
- [9] P. Ewen, A. Li, Y. Chen, S. Hong, and R. Vasudevan, "These Maps are Made for Walking: Real-Time Terrain Property Estimation for Mobile Robots," *IEEE Robot. Autom. Letter.*, vol. 7, no. 3, pp. 7083–7090, 2022.
- [10] A. Rosinol, M. Abate, Y. Chang, and L. Carlone, "Kimera: an Open-Source Library for Real-Time Metric-Semantic Localization and Mapping," in *Proc. IEEE Int. Conf. Robot. and Automation*, 2020, pp. 1689–1696.
- [11] S. T. O'Callaghan and F. T. Ramos, "Gaussian process occupancy maps," *Int. J. Robot. Res.*, vol. 31, no. 1, pp. 42–62, 2012.
- [12] K. Doherty, T. Shan, J. Wang, and B. Englot, "Learning-Aided 3-D Occupancy Mapping with Bayesian Generalized Kernel Inference," *IEEE Trans. Robot.*, vol. 35, no. 4, pp. 953–966, 2019.
- [13] W. R. Vega-Brown, M. Doniec, and N. G. Roy, "Nonparametric Bayesian inference on multivariate exponential families," in *Proc. Advances Neural Inform. Process. Syst. Conf.*, Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K. Weinberger, Eds., vol. 27, 2014.
- [14] P. Wu, S. Chen, and D. N. Metaxas, "MotionNet: Joint Perception and Motion Prediction for Autonomous Driving Based on Bird's Eye View Maps," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2020, pp. 11 382–11 392.
- [15] A. Saha, O. Mendez, C. Russell, and R. Bowden, "Translating images into maps," in *Proc. IEEE Int. Conf. Robot. and Automation*, 2022, pp. 9200–9206.
- [16] V. Cartillier, Z. Ren, N. Jain, S. Lee, I. Essa, and D. Batra, "Semantic MapNet: Building Allocentric SemanticMaps and Representations from Egocentric Views," in *Proc. AAAI Nat. Conf. Artif. Intell.*, February 2021.
- [17] J. Wilson, J. Song, Y. Fu, A. Zhang, A. Capodieci, P. Jayakumar, K. Barton, and M. Ghaffari, "MotionSC: Data Set and Network for Real-Time Semantic Mapping in Dynamic Environments," *IEEE Robot. Autom. Letter.*, vol. 7, no. 3, pp. 8439–8446, 2022.
- [18] A. Unnikrishnan, J. Wilson, L. Gan, A. Capodieci, P. Jayakumar, K. Barton, and M. Ghaffari, "Dynamic semantic occupancy mapping using 3D scene flow and closed-form Bayesian inference," *IEEE Access*, vol. 10, pp. 97 954–97 970, 2022.

- [19] J. Wilson, Y. Fu, A. Zhang, J. Song, A. Capodiecici, P. Jayakumar, K. Barton, and M. Ghaffari, "Convolutional Bayesian Kernel Inference for 3D Semantic Mapping," in *Proc. IEEE Int. Conf. Robot. and Automation*, 2023, pp. 8364–8370.
- [20] P. Jiang, P. Osteen, M. Wigness, and S. Saripalli, "RELLIS-3D Dataset: Data, Benchmarks and Analysis," in *Proc. IEEE Int. Conf. Robot. and Automation*, 2021, p. 1110–1116.
- [21] J. Stücker, N. Biresev, and S. Behnke, "Semantic mapping using object-class segmentation of RGB-D images," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots and Syst.*, 2012, pp. 3005–3010.
- [22] H. He and B. Upcroft, "Nonparametric semantic segmentation for 3D street scenes," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots and Syst.*, 2013, pp. 3697–3703.
- [23] S. Sengupta, E. Greveson, A. Shahrokni, and P. H. S. Torr, "Urban 3D semantic modelling using stereo vision," in *Proc. IEEE Int. Conf. Robot. and Automation*, 2013, pp. 580–585.
- [24] S. Sengupta and P. Sturgess, "Semantic octree: Unifying recognition, reconstruction and representation via an octree constrained higher order MRF," in *Proc. IEEE Int. Conf. Robot. and Automation*, 2015, pp. 1874–1879.
- [25] Z. Zhao and X. Chen, "Building 3D semantic maps for mobile robots using RGB-D camera," *Intell. Service Robot.*, vol. 9, 10 2016.
- [26] A. Kundu, Y. Li, F. Dellaert, F. Li, and J. M. Rehg, "Joint Semantic Segmentation and 3D Reconstruction from Monocular Video," in *Proc. European Conf. Comput. Vis.*, 2014, pp. 703–718.
- [27] J. Wang and B. Englot, "Fast, accurate gaussian process occupancy maps via test-data octrees and nested Bayesian fusion," in *Proc. IEEE Int. Conf. Robot. and Automation*, 2016, pp. 1003–1010.
- [28] M. G. Jadidi, L. Gan, S. A. Parkison, J. Li, and R. M. Eustice, "Gaussian Processes Semantic Map Representation," *ArXiv*, vol. abs/1707.01532, 2017.
- [29] N. Hughes, Y. Chang, and L. Carlone, "Hydra: A Real-time Spatial Perception System for 3D Scene Graph Construction and Optimization," in *Proc. Robot.: Sci. Syst. Conf.*, 06 2022.
- [30] H. Oleynikova, Z. Taylor, M. Fehr, R. Siegwart, and J. Nieto, "Voxblox: Incremental 3D Euclidean Signed Distance Fields for on-board MAV planning," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots and Syst.*, 2017, pp. 1366–1373.
- [31] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. u. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds., vol. 30. Curran Associates, Inc., 2017.
- [32] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Minneapolis, Minnesota: Association for Computational Linguistics, Jun. 2019, pp. 4171–4186.
- [33] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale," in *Proc. Int. Conf. Learning Representations*, 2021.
- [34] X. Zhu, W. Su, L. Lu, B. Li, X. Wang, and J. Dai, "Deformable {detr}: Deformable transformers for end-to-end object detection," in *Proc. Int. Conf. Learning Representations*, 2021.
- [35] L. Sun, Z. Yan, A. Zaganidis, C. Zhao, and T. Duckett, "Recurrent-OctoMap: Learning State-Based Map Refinement for Long-Term Semantic Mapping with 3-D-Lidar Data," *IEEE Robot. Autom. Letter.*, vol. 3, no. 4, pp. 3749–3756, 2018.
- [36] Y. Xiang and D. Fox, "DA-RNN: Semantic Mapping with Data Associated Recurrent Neural Networks," in *Robotics. Sci. Sys.*, vol. 13, 2017.
- [37] P. Sun, H. Kretzschmar, X. Dotiwalla, A. Chouard, V. Patnaik, P. Tsui, J. Guo, Y. Zhou, Y. Chai, B. Caine, V. Vasudevan, W. Han, J. Ngiam, H. Zhao, A. Timofeev, S. Ettinger, M. Krivokon, A. Gao, A. Joshi, Y. Zhang, J. Shlens, Z. Chen, and D. Anguelov, "Scalability in Perception for Autonomous Driving: Waymo Open Dataset," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, June 2020.
- [38] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liong, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom, "nuScenes: A Multimodal Dataset for Autonomous Driving," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2020, pp. 11 618–11 628.
- [39] 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 *Proc. IEEE/RSJ Int. Conf. Intell. Robots and Syst.*, 2019, pp. 5000–5007.
- [40] D. Maturana, P.-W. Chou, M. Uenoyama, and S. Scherer, "Real-Time Semantic Mapping for Autonomous Off-Road Navigation," in *Field and Serv. Robot.*, M. Hutter and R. Siegwart, Eds. Cham: Springer International Publishing, 2018, pp. 335–350.
- [41] H. Thomas, C. R. Qi, J.-E. Deschaud, B. Marcotegui, F. Goulette, and L. Guibas, "KPConv: Flexible and Deformable Convolution for Point Clouds," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2019, pp. 6410–6419.
- [42] T. Cortinhal, G. Tzelepis, and E. Erdal Aksoy, "SalsaNext: Fast, Uncertainty-Aware Semantic Segmentation of LiDAR Point Clouds," in *Advances in Visual Computing*, 2020, pp. 207–222.
- [43] L. D. Jackel, E. Krotkov, M. Perschbacher, J. Pippine, and C. Sullivan, "The DARPA LAGR program: Goals, challenges, methodology, and phase I results," *J. Field Robot.*, vol. 23, 2006.
- [44] A. Kelly, A. Stentz, O. Amidi, M. Bode, D. Bradley, A. Diaz-Calderon, M. Happold, H. Herman, R. Mandelbaum, T. Pilarski, P. Rander, S. Thayer, N. Vallidis, and R. Warner, "Toward Reliable Off Road Autonomous Vehicles Operating in Challenging Environments," *Int. J. Robot. Res.*, vol. 25, no. 5-6, pp. 449–483, 2006.
- [45] R. Manduchi, A. Castano, A. Talukder, and L. H. Matthies, "Obstacle Detection and Terrain Classification for Autonomous Off-Road Navigation," *Auton. Robot.*, vol. 18, pp. 81–102, 2005.
- [46] A. Melkumyan and F. Ramos, "A Sparse Covariance Function for Exact Gaussian Process Inference in Large Datasets," in *Proc. Int. Joint Conf. Artif. Intell.*, 2009, p. 1936–1942.
- [47] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications," *ArXiv*, vol. abs/1704.04861, 2017.
- [48] E. Gordon-Rodriguez, G. Loaiza-Ganem, and J. P. Cunningham, "The continuous categorical: a novel simplex-valued exponential family," in *Proceedings of the 37th International Conference on Machine Learning*, ser. ICML'20, 2020.
- [49] J. Kim, J. Seo, and J. Min, "Evidential Semantic Mapping in Off-road Environments with Uncertainty-aware Bayesian Kernel Inference," *arXiv*, vol. abs/2403.14138, 2024.
- [50] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "PointNet++: Deep Hierarchical Feature Learning on Point Sets in a Metric Space," in *Proc. Advances Neural Inform. Process. Syst. Conf.*, 2017, pp. 1–10.
- [51] Y. Wang, Y. Sun, Z. Liu, S. E. Sarma, M. M. Bronstein, and J. M. Solomon, "Dynamic Graph CNN for Learning on Point Clouds," *IEEE Trans. Graph.*, vol. 38, no. 5, oct 2019.
- [52] C. Rasmussen and C. Williams, *Gaussian processes for machine learning*. MIT press, 2006, vol. 1.
- [53] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, "PyTorch: An Imperative Style, High-Performance Deep Learning Library," in *Proc. Advances Neural Inform. Process. Syst. Conf.*, 2019, pp. 8024–8035.
- [54] R. A. Newcombe, S. Izadi, O. Hilliges, D. Molyneaux, D. Kim, A. J. Davison, P. Kohi, J. Shotton, S. Hodges, and A. Fitzgibbon, "KinectFusion: Real-time dense surface mapping and tracking," in *Proc. IEEE Int. Symp. Mixed and Augm. Reality*, 2011, pp. 127–136.
- [55] T. Ort, J. M. Walls, S. A. Parkison, I. Gilitschenski, and D. Rus, "Maplite 2.0: Online hd map inference using a prior sd map," *IEEE Robot. Autom. Letter.*, vol. 7, no. 3, pp. 8355–8362, 2022.
- [56] T. Shan, B. Englot, D. Meyers, W. Wang, C. Ratti, and R. Daniela, "LIO-SAM: Tightly-coupled Lidar Inertial Odometry via Smoothing and Mapping," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots and Syst.*, 2020, pp. 5135–5142.
- [57] P. Fankhauser and M. Hutter, "A Universal Grid Map Library: Implementation and Use Case for Rough Terrain Navigation," in *Robot Operating System (ROS) – The Complete Reference*, 2016, vol. 1, ch. 5.
- [58] D. P. Kingma and J. Ba, "Adam: A Method for Stochastic Optimization," in *Proc. Int. Conf. Learning Representations*, Y. Bengio and Y. LeCun, Eds., 2015.
- [59] S. Sengupta, E. Greveson, A. Shahrokni, and P. H. S. Torr, "Urban 3D semantic modelling using stereo vision," in *Proc. IEEE Int. Conf. Robot. and Automation*, 2013, pp. 580–585.
- [60] F. Yu and V. Koltun, "Multi-Scale Context Aggregation by Dilated Convolutions," in *Proc. Int. Conf. Learning Representations*, 2016.
- [61] S. Yang, Y. Huang, and S. Scherer, "Semantic 3D occupancy mapping through efficient high order CRFs," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots and Syst.*, 09 2017, pp. 590–597.

- [62] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for Autonomous Driving? The KITTI Vision Benchmark Suite," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2012, pp. 3354–3361.
- [63] A. Geiger, M. Roser, and R. Urtasun, "Efficient large-scale stereo matching," in *Proc. Asian Conf. Comput. Vis.*, 2011, pp. 25–38.
- [64] R. Mur-Artal, J. Montiel, and J. Tardos, "Orb-slam: a versatile and accurate monocular slam system," *IEEE Trans. Robot.*, vol. 31, pp. 1147 – 1163, 10 2015.
- [65] K. Cho, B. van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, "Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation," in *Proc. Conf. Empirical Methods in NLP*, 2014, pp. 1724–1734.
- [66] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, p. 1735–1780, 1997.
- [67] H. Tang, Z. Liu, S. Zhao, Y. Lin, J. Lin, H. Wang, and S. Han, "Searching Efficient 3D Architectures with Sparse Point-Voxel Convolution," in *Proc. European Conf. Comput. Vis.*, 2020.



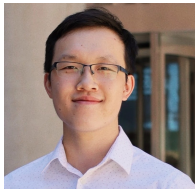
Andrew Capodici is the Director of Robotics at Neya Systems and has been with Neya for over a decade. During this time, Andrew has performed and managed applied research in all areas of the robotics stack including kinodynamically-feasible path planning in congested spaces, traversability estimation in off-road terrain, and negative obstacle detection. As Director of Robotics, Andrew is focused on transitioning Neya's state-of-the-art autonomy research into fieldable, robust autonomy capabilities that deliver value to the warfighter and commercial off-road spaces. Andrew has led numerous multi-million-dollar programs including Neya's work on GVSC's Combat Vehicle Robotics program, and programs to develop autonomous construction vehicles for the commercial sector.



Joey Wilson received the B.S. degree in computer engineering from California Polytechnic State University San Luis Obispo (Cal Poly), CA, USA, in 2019. He is currently a Ph.D. candidate in the University of Michigan Robotics Department, Ann Arbor, MI, USA. His research interests include spatial memory in dynamic and uncertain environments for autonomous systems.



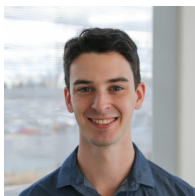
Paramsothy Jayakumar received the B.Sc.Eng. degree (Hons.) from the University of Peradeniya, Sri Lanka, and the M.S. and Ph.D. degrees from Caltech. He worked at Ford Motor Company and BAE Systems. He is a Senior Technical Expert of analytics with U.S. Army DEVCOM Ground Vehicle Systems Center (GVSC). He has published over 200 papers in peer-reviewed literatures. He is a fellow of the Society of Automotive Engineers and the American Society of Mechanical Engineers. He received the DoD Laboratory Scientist of the Quarter Award, the NATO Applied Vehicle Technology Panel Excellence Awards, the SAE Arch T. Colwell Cooperative Engineering Medal, the SAE James M. Crawford Technical Standards Board Outstanding Achievement Award, the BAE Systems Chairman's Award, and the NDIA GVSETS Best Paper Awards. He is also an Associate Editor of the ASME Journal of Autonomous Vehicles and Systems, and the Editorial Board Member of the International Journal of Vehicle Performance and the Journal of Terramechanics.



Yuewei Fu received the B.S. degree in mechanical engineering from New York University (NYU) in 2021. He is currently a M.S. student in the University of Michigan Robotics Department, Ann Arbor, MI, USA. His research interests include scene understanding for off-road and underwater robots in dynamic environments.



Kira Barton received the Ph.D. degree in Mechanical Engineering from the University of Illinois at Urbana-Champaign, USA, in 2010. She is currently a Professor at the Robotics Institute and Department of Mechanical Engineering, University of Michigan, Ann Arbor, MI, USA. Her research interests lie in control theory and applications including high precision motion control, iterative learning control, and control for autonomous vehicles.



Josh Friesen received the B.S. degree in software engineering from Fresno Pacific University in 2021. He is currently a M.S. student in the University of Michigan Robotics Department, Ann Arbor, MI, USA. His research interests include deep learning for autonomous robots.



Maani Ghaffari received the Ph.D. degree from the Centre for Autonomous Systems (CAS), University of Technology Sydney, NSW, Australia, in 2017. He is currently an Assistant Professor at the Department of Naval Architecture and Marine Engineering and the Department of Robotics, University of Michigan, Ann Arbor, MI, USA. He is the director of the Computational Autonomy and Robotics Laboratory. He is the recipient of the 2021 Amazon Research Awards. His research interests lie in the theory and applications of robotics and autonomous systems.



Parker Ewen received his MSc in Robotics, Systems, and Control from the ETH Zurich in 2020. He is currently a PhD candidate at the University of Michigan where he is a member of ROAHM Lab. He researches numerical techniques for state estimation, mapping, and active learning with applications for robotics.