

WildFusion: Multimodal Implicit 3D Reconstructions in the Wild

Yanbaihui Liu¹ and Boyuan Chen¹
www.generalroboticslab.com/WildFusion

Abstract—We propose WildFusion, a novel approach for 3D scene reconstruction in unstructured, in-the-wild environments using multimodal implicit neural representations. WildFusion integrates signals from LiDAR, RGB camera, contact microphones, tactile sensors, and IMU. This multimodal fusion generates comprehensive, continuous environmental representations, including pixel-level geometry, color, semantics, and traversability. Through real-world experiments on legged robot navigation in challenging forest environments, WildFusion demonstrates improved route selection by accurately predicting traversability. Our results highlight its potential to advance robotic navigation and 3D mapping in complex outdoor terrains.

I. INTRODUCTION

Robots need effective environmental representations to navigate safely and accomplish tasks successfully in unstructured outdoor environments – often referred to as “in-the-wild” settings such as monitoring high-voltage power lines and extinguishing forest fires. However, accurately modeling these environments presents significant challenges due to their inherent complexity. The lack of clear boundaries between objects, combined with fluctuating lighting conditions and shadows, creates a dynamic and often ambiguous visual landscape. Additionally, these environments don’t follow predefined patterns or rules for interpretation, hence creating models that can reliably understand and represent them becomes difficult.

Traditional 3D reconstruction methods, such as LOAM [1] and LIO-SAM [2], typically rely on a single sensor modality like LiDAR or cameras. These methods have shown success in mapping static scenes. Other approaches [3–5] employ semantic segmentation to detect and handle moving objects for improved performance in dynamic environments. However, both approaches struggle to provide high-quality and drift-free reconstruction in complex in-the-wild settings, where a single sensor modality is not enough to overcome the scene complexities and the inherent limitations of vision-based sensors.

Recent research has focused on addressing this issue by studying multimodal sensor fusion. This includes combining LiDAR and cameras [6, 7], or integrating additional modalities such as audio, language and, tactile data [8–10]. While these multi-sensing systems allow for scene reconstructions for more advanced tasks in complex environments, they still rely on explicit representations, such as point clouds [11], voxels [12], or meshes [13]. Techniques such as Gaussian Splatting [14, 15] have further advanced explicit methods for processing dynamic outdoor scenes with stronger robustness

^{*}This work is supported by DARPA TIAMAT HR00112490419, DARPA FoundSci HR00112490372, ARL STRONG W911NF2320182 and W911NF2220113. ¹ All authors are from Duke University.

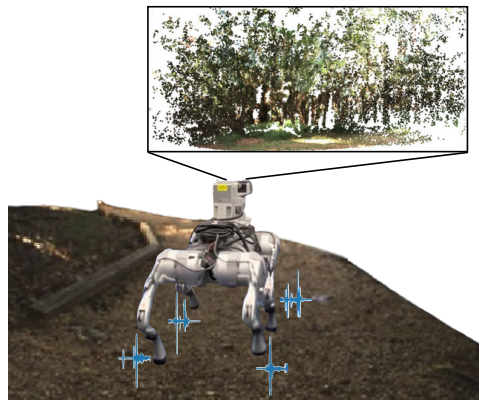


Fig. 1: **WildFusion** integrates LiDAR, camera, microphones, and tactile sensors with implicit neural representations for continuous 3D scene reconstruction. The learned model provides accurate and continuous traversability predictions to enhance legged robot navigation in forest environments.

to noise and more effective handling sparse data. However, as other explicit representations, they still require dense data to produce accurate models.

In recent years, implicit representation methods using Neural Radiance Fields (NeRF) [16–19] and Signed Distance Fields (SDFs) [20, 21] have received increasing attention. These approaches are particularly useful for describing complex surfaces and are more robust to sparse or incomplete data. Although existing implicit methods have made breakthroughs in 3D modeling and mapping, most solutions still rely heavily on visual or geometric sensor data. Research on fusing multimodal signals, such as acoustic and tactile data, remains limited, even though these modalities are highly useful for providing richer sensory information in unstructured environments.

We introduce WildFusion (Fig. 1), a multimodal implicit 3D reconstruction framework that integrates LiDAR, RGB camera, audio, and tactile sensors on a quadruped robot platform. WildFusion generates pixel-level continuous scene representations of multiple environmental features, including 3D geometry, semantics, continuous traversability, color, and confidence scores. By leveraging implicit representations, our method can produce complete scene representations from sparse inputs. We present a multimodal in-the-wild dataset by navigating our quadruped robot in a forest environment. Our experiments suggest that our multimodal formulation effectively leverages all modalities to enhance the 3D scene reconstruction with a richer understanding of the surrounding environment. By testing WildFusion on a downstream legged robot navigation task, we show that the robot can safely tra-

verse various terrains such as grasslands, dense leaves, high vegetation, and gravel. Furthermore, since our model design is modular, we believe that WildFusion provides a promising direction to integrate more diverse sensing modalities to build rich and robust scene representations.

II. RELATED WORK

A. 3D Mapping

Simultaneous Localization and Mapping (SLAM) is a key technique for mapping and navigating in unknown environments. SLAM enables robots to build maps of the environment while determining their positions. Traditional SLAM algorithms, such as direct methods [22] and feature-based methods [23, 24], have been successfully applied to city-view mapping. However, these methods often rely on salient environmental features which may not be available in unstructured environments. Moreover, they also produce sparse maps without sufficient details for more complex environments and tasks. Other works focus on generating dense maps [25, 26] to incorporate scene geometry and color to provide richer environmental information, but still leave gap and blank areas. Moreover, the use of explicit representations makes the system sensitive to outliers.

To address the above challenges, researchers have begun exploring implicit representations to enhance SLAM. Some approaches use radiance-based or coordinate-based representations [27–30] to effectively fill surface holes and represent both geometry and appearance such as color and texture in a photorealistic way. Others rely on SDF-based methods [31–34] for fine-grained surface reconstruction. Our approach builds on these advancements by leveraging implicit representations in SLAM to achieve continuous environmental mapping. However, we go further by incorporating multimodal signals beyond just vision with real-world legged robot demonstrations.

B. Traversability Prediction

Traversability prediction in outdoor settings has been addressed through several methods. The most straightforward and efficient approaches directly derive traversability from geometric representations or elevation maps [35–37]. However, these techniques often struggle when faced with similar terrains in the same area. Additionally, a heavy reliance on geometry ignores textures and semantics, which leads to mistakes hindering smooth navigation in outdoor environments, such as classifying high vegetation as solid obstacles. Semantic-aided methods help mitigate these issues by assigning different traversability costs to distinct terrain classes [38–40]. However, the predicted traversability values cannot handle fine-grained classes such as wet soil or dry land – two surfaces that belong to the same high-level terrain classes but offer very different levels of ease for traversal.

More advanced methods predict traversability by considering the robot’s physical interactions with the environment. For example, contrastive learning with human driving data and segmentation masks [41] has been used to predict

traversability in both on- and off-trail settings. Terrain semantics can also be mapped to vehicle speed profiles [42], allowing robots to adapt their speed in real time as terrain conditions change. Other approaches integrate vision and confidence estimation to account for deviations between expected and actual speeds, improving navigation in uncertain environments [43]. Additionally, risk-aware frameworks penalize uncertain terrains using evidential learning to enhance decision-making by factoring in traction predictions [44]. Our method also leverages the physical interactions with the environment from a legged robot, but instead uses proprioception signals to compute ground-truth continuous traversability scores. Our design enables us to generate pixel-level traversability predictions that accurately reflect the robot’s interaction with the environment.

C. Acoustic Profiling

Sound has emerged as a valuable modality for object perception. Previous research has shown that feedback from audio synthesis engines can be used to predict physical parameters of objects [45], and acoustic vibrations can be used to sense real-world object properties [46, 47]. In terrain classification, sound has been successfully integrated into robotic systems. For instance, bat-inspired echolocation has been used to accurately classify terrain types [48], such as grass, concrete, sand, and gravel, by using signal filtering and machine learning techniques. Similarly, acoustic sensors have been applied to real-world terrain classification tasks in robotics [49, 50], and multimodal approaches combining sound data with cameras and foot sensors have improved semantic predictions [49, 51]. Despite its potential, sound has not been widely used to support traversability prediction. To address this gap, we incorporate acoustic vibrations into our mapping system to predict environmental conditions and the understanding of semantic information of terrains that are in contact with the robot’s body.

III. METHOD

A. System Overview

WildFusion aims to represent in-the-wild environments by constructing a multimodal 3D map. Our approach integrates multiple sensor groups as input and output modalities, including contact microphones, camera, LiDAR, tactile, and proprioception sensors. By exploring synergistic relationships between these modalities, our model gains a rich and accurate understanding of the complex environment. By combining these diverse perceptual signals, WildFusion addresses the limitations of traditional vision-based only methods and offers a more robust framework for scene representation and navigation.

Our system (Fig. 2) processes colored point clouds by fusing RGB images with LiDAR data, and Mel spectrograms from contact microphones, and random sample coordinates as query points. The system outputs semantic category, color, traversability, SDF, and confidence scores for each queried point. Intuitively, our system answers multiple environmental characteristics from multimodal perceptual inputs for the

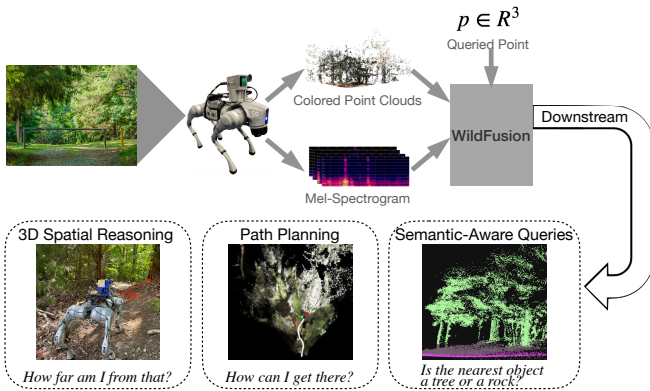


Fig. 2: **System Overview:** WildFusion leverages signals from a RGB camera, LiDAR, contact microphones, tactile sensors, and an IMU to build rich multimodal scene representations. The learned representation enables the creation of dense maps with diverse features and guides legged robot navigation through continuous traversability scores.

queried points. Therefore, WildFusion can be interpreted as a queried-based differentiable scene representation.

B. Integrated Multimodal Robotic Platform

We built a legged robotic platform (Fig. 3) capable of multimodal perception based on Unitree Go2. In addition to the original IMU and on-foot tactile sensors, we equipped the platform with several key modules. We first added a Blackfly S USB3 RGB camera due to its flexibility for lens mount, allowing us to optimize the field of view and depth of field with interchangeable lens. We paired it with a Tamron C-Mount 4 to 12mm Varifocal Manual Iris lens, which provides adjustable focal length and iris control for precise image capture tailored to our robot platform and forest setting. Additionally, we mounted a Livox AVIA solid-state LiDAR next to the camera with a 3D-printed frame. Our selected LiDAR sensor has a unique non-repetitive scanning pattern, which captures richer data over time compared to traditional mechanical LiDAR systems. This is especially useful in complex forest scenes, as it generates dense point clouds that improve the accuracy of ground-truth 3D information such as SDF. The sensor is also lightweight and cost-effective, offering higher field-of-view coverage within 0.3s at just a fraction of the cost of traditional 64-line LiDARs.

As the robot moves across different terrains, contact between its feet and the ground will generate distinct vibrations. To capture these physical interactions without interference from environmental noises, we chose piezoelectric contact microphones and mounted them on the robot’s calves. All sensors are connected to the onboard computer with synchronized timestamps. To ensure stable locomotion on challenging forest terrains, we trained a walking policy using a reinforcement learning strategy [52]. The robot walks at a constant speed of 0.4m/s on a flat concrete floor, with variations in speed when traversing different terrain types.

C. WildFusion Dataset

The raw dataset includes RGB images, point clouds, audio recordings, foot contact forces, and IMU data collected from

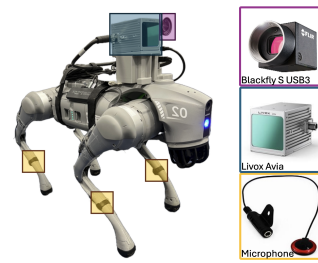


Fig. 3: **Our Robotic Platform:** We integrate a monocular RGB camera, LiDAR and contact microphones on the Go2 platform. Proprioceptive data from IMU and foot tactile sensors is also recorded to calculate ground truth for our model training.

the quadruped robot over a 750m trajectory during a 30min walk through a forest inside Eno River State Park in NC, USA. This scene covers 90% forested area with trees, leaves, tall vegetation, and fallen tree trunks, while the remaining 10% consists of gravel or concrete paths. Due to the non-repeated scanning pattern of the AVIA LiDAR, each frame is generated by fusing data collected over a 2-second interval. In total, the processed dataset contains 551 frames, with 460 used for training, 50 for validation, and 41 for testing. Ground truth labels are provided for semantic segmentation, traversability, SDF, and confidence scores for each frame.

Semantic Labels are manually assigned to points on or within objects, without differentiating between individual instances. **Color** is post-processed by converting it from RGB to LAB color space due to its uniformity of color perceptual distances and better handling of illumination variations. **Traversability** scores are calculated frame-by-frame using accelerometer and tactile sensor data. Accelerometer readings are normalized across all three axes, and their variance is used to assess movement instability. Tactile data is processed by normalizing the force distribution across four sensors and calculating its deviation from an ideal balanced state, which indicates stable ground contact from a previous calibration recording. The final traversability score for each frame is the product of the accelerometer variance and tactile deviation with a normalized continuous value between 0 and 1.

To compute the **SDF** from LiDAR data, points are uniformly sampled along each LiDAR ray, and the distance to the nearest surface point is calculated using a K-D Tree with 30 leaf nodes. Free space points are labeled based on their minimum distance to observed surface points, while negative distances are determined by sampling points beyond the termination of the LiDAR ray to account for the finite depth of objects. Confidence scores are derived from the SDF, where free space points are assigned a confidence score of 1. For points with negative SDF values, the confidence score decreases exponentially as the distance from the surface increases. Semantic labels and color data for free space points are set to NULL to reflect the absence of real-world information in these regions.

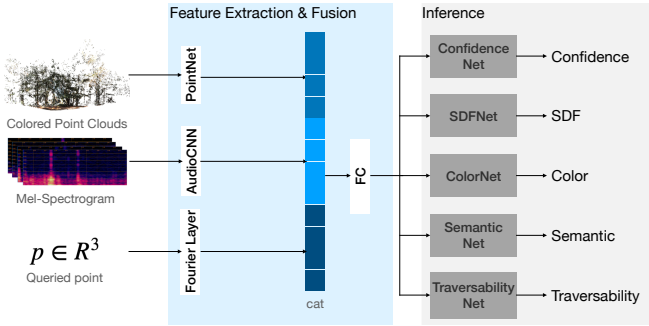


Fig. 4: **WildFusion Model:** The model takes in a colored point cloud, Mel spectrogram, and queried points through a combination of encoders. The final outputs are derived through specialized sub-networks to predict confidence, SDF, color, semantics, and traversability.

D. WildFusion Model

To overcome the limitations of traditional vision-based only methods for scene understanding in unstructured environments, WildFusion introduces a multimodal implicit scene representation learning framework. Our key idea is to not only condition multimodal signals from the input, but also ask the model to fuse these features to output multimodal environmental features. Furthermore, our model is established as implicit neural representations to enable dense predictions from sparse inputs by simply querying the model with coordinates from higher resolutions. The diverse data types complement one another, hence they enhance the map’s details and utility for downstream tasks. For instance, traditional approaches that only predict SDF are limited by ambiguous object boundaries and insufficient information to learn complex geometries. Instead, we train our model to simultaneously predict color and semantics along with SDF, so that these additional modalities help inform clearer object boundaries. Confidence scores give our model the ability to assess and express its uncertainty about SDF predictions. Additionally, traversability predictions derived from acoustic vibrations can assist in predicting the geometry and semantics of the ground near the robot.

We show our model architecture in Fig. 4. Specifically, the model takes colored point clouds, Mel-spectrograms, and queried coordinates as input. The queried point is passed through Fourier feature encoding [53] to efficiently capture high-frequency functions in low-dimensional data. Point clouds, consisting of 15,000 points, are processed using a modified PointNet [54]. We incorporate a T-Net for both input and feature transformations to ensure rotation invariance. We reduce the feature dimensions from 1024 to 512 and introduce residual connections to enhance gradient flow and convergence. The point cloud features capture both global and hierarchical information about the environment. Acoustic vibrations are converted to Mel-spectrograms, using a Fast Fourier Transform (FFT) window size of 2048 and 128 Mel bands, with a hop length of 512 and the highest frequency limited to 8,192Hz. These spectrograms, generated from 0.5-second audio segments, are then stacked along the

TABLE I: Quantitative Results

		Seen Scenes (w/ diff. viewpoints)	Unseen Scenes
Color	MSE	0.056	0.064
	MAE	0.149	0.167
	PSNR	13.690	12.258
Geometry	Haus.	5.621	6.288
	Cha.	0.072	0.079
Semantic	Acc.	0.886	0.755
	Pre.	0.924	0.808
	Re.	0.891	0.755
	F1	0.907	0.726
	IoU	0.825	0.692
confidence	ECE	0.189	0.217

leg channel and passed through several convolutional layers to extract relevant audio features.

Features extracted from different modalities are concatenated and fed into a comprehensive model layer, which predicts SDF, confidence score, color, traversability, and semantic through separate branches. The SDF and confidence score predictions utilize fully connected networks with residual skip connections to maintain smooth training dynamics. The SDF prediction applies a Tanh activation for smooth and bounded outputs, while the confidence score prediction uses a Sigmoid activation to map values between 0 and 1. For color and semantic label predictions, fully connected layers with ReLU are used to encourage efficient learning. Minimal dropout is applied in these branches to prevent overfitting while ensuring the model generalizes well.

We train our model using end-to-end supervised learning on our dataset. We minimize the following loss function which is a weighted sum of the losses from each branch:

$$L = \lambda_1 L_{\text{SDF}} + \lambda_2 L_{\text{eikonal}} + \lambda_3 L_{\text{confidence}} + \lambda_4 L_{\text{semantics}} + \lambda_5 L_{\text{color}} + \lambda_6 L_{\text{traversability}}$$

where λ values are constants that balance different losses. L_{SDF} is Huber loss and we also include a small Eikonal loss [55] to ensure the gradient norm of SDF is close to one. This encourages the model to learn a geometrically consistent implicit surface representation. $L_{\text{confidence}}$ helps focus on uncertain regions by assigning different weights based on whether the target confidence equals one:

$$\frac{1}{N} \sum_{i=1}^N [\alpha \cdot \mathbb{I}(c_i = 1) \cdot (\hat{c}_i - c_i)^2 + \beta \cdot \mathbb{I}(c_i \neq 1) \cdot (\hat{c}_i - c_i)^2]$$

where \hat{c}_i is the predicted confidence for the i^{th} sample, c_i is the ground truth confidence for the i^{th} sample, $\mathbb{I}(\cdot)$ is an indicator function, and α and β are weights corresponding to whether the confidence is 1 or not. $L_{\text{semantics}}$ and L_{color} are cross entropy loss, and the $L_{\text{traversability}}$ uses mean squared error loss. Our model was trained on $8 \times$ NVIDIA A6000 GPUs for about 9 hours. We will open source all of our dataset, hardware integration solution, and code.

IV. EXPERIMENT

A. Geometry, Confidence, Color, and Semantics

To evaluate the performance of our model, we conducted experiments on two groups of test sets: previously visited

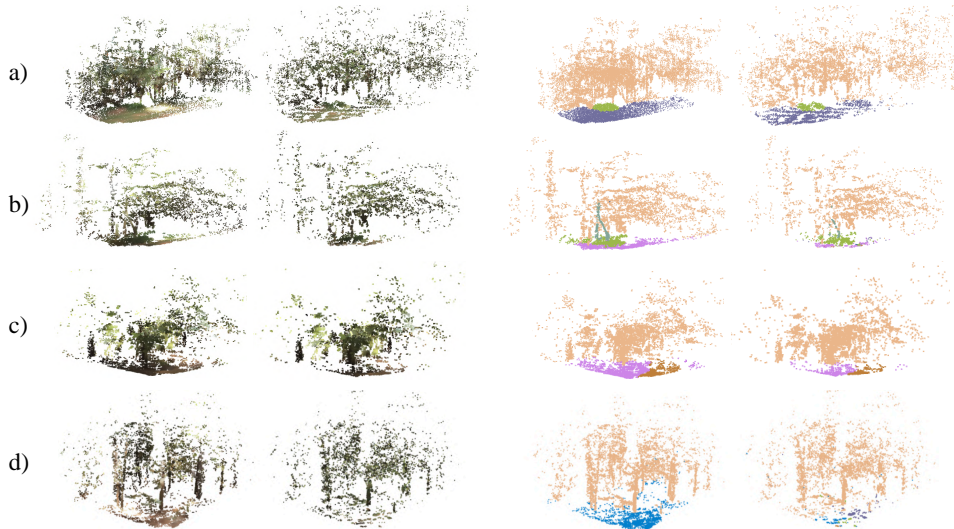


Fig. 5: Visualization comparing geometry, color, and semantic predictions: The column 1 and 3 display the ground truth geometry, color, and semantics, while the column 2 and 4 show the corresponding predictions. a) and b) represent previously visited location from a different viewpoint. c) and d) represent unvisited location during training.

locations but observed from slightly different viewpoints and entirely unvisited locations. The previously visited locations with novel views assess the prediction consistency, while the unvisited locations test the model’s generalizability and robustness. We sampled 30,000 points for each frame to get the prediction result. To ensure the validity of our results, we focused only on sample points with SDF values less than or equal to zero, as well as points with valid predicted semantic IDs. These criteria indicate that the points are either on the surface or inside objects, ensuring they provide meaningful data for geometry and semantic evaluation.

We employed a range of evaluation metrics across four key dimensions: color, geometry, confidence, and semantics. For color predictions, we report the mean-squared error (MSE), mean absolute error (MAE), and peak signal-to-noise ratio (PSNR). For geometry predictions, we calculated the Hausdorff distance to measure maximum point-based deviations and captured worst-case scenarios. We also computed the Chamfer distance to get average nearest-point distances. For semantic predictions, we applied standard classification metrics including accuracy, precision, recall, F1 score, and intersection over union (IoU). For confidence, we use expected calibration error (ECE).

As presented in Table I, our model achieved high accuracy in predicting geometry, confidence, color, and semantics for previously visited locations with unseen viewpoints. For unseen scenes, while some discrepancies arose compared to the ground truth, the model still produced meaningful predictions. Fig. 5 visualizes these results, indicating that although performance slightly decreases in unseen environments, the model maintained a reasonable level of accuracy and successfully generalized to new scenes.

B. Modality Studies

Our hypothesis is that all modalities help the model learn richer and more accurate 3D representations. We ablated our

TABLE II: Modality Studies

		Full mode	SDF, Semantic, Color	SDF, Conf. Color	SDF, Conf., semantic	SDF, Conf.
Color	MSE	0.060	0.061	0.066	-	-
	MAE	0.162	0.172	0.176	-	-
	PSNR	12.633	12.518	12.097	-	-
Geometry	Haus.	6.166	6.499	6.512	6.512	7.514
	Cha.	0.068	0.081	0.082	0.083	0.083
	Acc.	0.769	0.717	-	0.705	-
Semantic	Prec.	0.851	0.807	-	0.774	-
	Re.	0.763	0.717	-	0.705	-
	F1	0.770	0.694	-	0.693	-
	IoU	0.704	0.638	-	0.619	-
Confidence	ECE	0.205	-	0.79	0.228	0.778

design by removing various modalities and compared the results with our original model (Table II). Our results show that all modalities contribute to the overall performance. The absence of confidence scores introduced additional uncertainty into predictions, leading to less accurate predictions. Similarly, the removal of color or semantics leads to a drop in shape reconstruction accuracy, as evidenced by increased Chamfer and Hausdorff distances. These two modalities show an interdependence, as both are crucial for capturing fine-grained details. Semantics, in particular, helped the model become more confident in its predictions. The consistently better performance of our full model, compared to the ablated versions, confirms the complementary nature of each modality in producing richer and more accurate 3D representations.

C. Motion Planning and Traversability Analysis

Our model can be effectively used for downstream motion planning tasks, as we will demonstrate on our legged robot. Since the direct output from our model provides a single traversability score per scene, we refine this score to obtain pixel-level scores for motion planning by grounding it in scene context. First, we leverage the semantics of objects in the scene. Each semantic class in our dataset is manually assigned a traversability score, allowing us to generate a semantics-based traversability mask based on our predicted

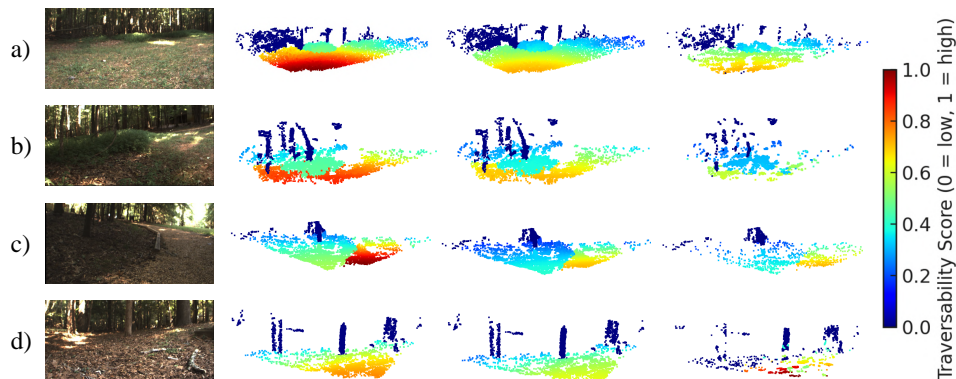


Fig. 6: Traversability: Columns from left to right show the real scene, weight factor, ground truth traversability, and predicted traversability. The weight factor scales the single traversability score into point-wise values based on both semantic labels and distance from the robot. a) and b) represent previously visited location from a different viewpoint. c) and d) represent unvisited location during training. Blue indicates non-traversable areas, and red indicates fully traversable areas.



Fig. 7: Real-world path planning: Compared the trajectory of our proposed model (solid robot) against baseline methods (60% transparent robot). The red dot indicates the desired destination, and yellow cross mark means the robot stuck here. a) shows comparison against elevation-based costmap. b) shows comparison against semantic-based costmap.

semantics. Second, because acoustics vibrations from foot-ground interactions influence our traversability prediction, we give higher weight to prediction near the robot’s location. This results in a distance-based traversability mask. Specifically, this mask is created by taking the Hadamard product of the distance-based mask and a Gaussian distribution matrix centered around the robot. We empirically determined the variance to be six since the variance decides how quickly the weights decrease. The final pixel-level traversability scores are calculated as the Hadamard product of the semantic and distance-based masks, weighted by our model’s prediction.

Fig. 6 provides a visualization of weight factors and prediction results. The overall performance of traversability prediction aligned well with the ground truth, though specific values deviated by approximately 5%. This discrepancy is likely due to acoustic vibration data can only inform more accurate traversability of the near-feet grounds. One potential solution is to study longer dependencies from the acoustic data with more sophisticated network designs in the future.

We conducted a real-world experiment to evaluate how traversability prediction guides robot motion planning. We used an A^* planner and projected our traversability scores

as the costmap. Regions with lower traversability score has higher cost and vice versa. We compared our model with two baselines: an elevation map-based costmap and a variant of our proposed model that only predicts semantic information and assigns cost based on category. Fig. 7 shows the performance of our method against baselines in multiple scenarios. The elevation map based method showed the most limitations. It misidentified dense vegetation as obstacles and stopped moving forward. Moreover, it failed to differentiate between terrains with similar height, thus won’t avoid the harder terrains. Our model performed well in both cases and would try to choose the safest route. Comparative analysis with the semantic-based method indicated our complete version of WildFusion enhanced the understanding of terrains by effectively guiding the robot toward flatter routes. In ambiguous scenarios, such as those combining solid ground with terrain littered with falling branches or grass, since the variation is different between them, our model’s ability to discern distinct vibration patterns allowed better decision-making.

V. CONCLUSIONS, LIMITATIONS, AND FUTURE WORK

WildFusion presents a novel multimodal framework for implicit 3D reconstruction and navigation in unstructured outdoor environments. By integrating data from LiDAR, RGB camera, contact microphone, and proprioceptive sensors, WildFusion creates a rich pixel-level representation of complex terrains, which enhances environmental understanding and provides support for more efficient and reliable path planning. Our experiments show that this approach significantly improves navigation in a forest setting.

For future work, expanding sensor modalities, such as incorporating humidity and thermal sensors, could provide a more holistic environmental understanding. Furthermore, refining our traversability prediction with more advanced models could lead to even greater accuracy by considering more dynamic information from the robot motions. Moreover, our current model requires offline computational resources. A future direction can study how to enable on-board training with edge computing techniques.

REFERENCES

- [1] J. Zhang and S. Singh, “Low-drift and real-time lidar odometry and mapping,” *Autonomous Robots*, vol. 41, no. 2, pp. 401–416, February 2017.
- [2] T. Shan, B. Englot, D. Meyers, W. Wang, C. Ratti, and D. Rus, “Lio-sam: Tightly-coupled lidar inertial odometry via smoothing and mapping,” in *2020 IEEE/RSJ international conference on intelligent robots and systems (IROS)*. IEEE, 2020, pp. 5135–5142.
- [3] B. Bescos, J. M. Fácil, J. Civera, and J. Neira, “Dyna-slam: Tracking, mapping, and inpainting in dynamic scenes,” *IEEE Robotics and Automation Letters*, vol. 3, no. 4, pp. 4076–4083, 2018.
- [4] C. Yu, Z. Liu, X.-J. Liu, F. Xie, Y. Yang, Q. Wei, and Q. Fei, “Ds-slam: A semantic visual slam towards dynamic environments,” in *2018 IEEE/RSJ international conference on intelligent robots and systems (IROS)*. IEEE, 2018, pp. 1168–1174.
- [5] X. Chen, A. Milioto, E. Palazzolo, P. Giguere, J. Behley, and C. Stachniss, “Suma++: Efficient lidar-based semantic slam,” in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2019, pp. 4530–4537.
- [6] J. Lin and F. Zhang, “R 3 live: A robust, real-time, rgb-colored, lidar-inertial-visual tightly-coupled state estimation and mapping package,” in *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022, pp. 10 672–10 678.
- [7] L. Kong, X. Xu, J. Ren, W. Zhang, L. Pan, K. Chen, W. T. Ooi, and Z. Liu, “Multi-modal data-efficient 3d scene understanding for autonomous driving,” *arXiv preprint arXiv:2405.05258*, 2024.
- [8] Y. Yue, C. Yang, J. Zhang, M. Wen, Z. Wu, H. Zhang, and D. Wang, “Day and night collaborative dynamic mapping in unstructured environment based on multimodal sensors,” in *2020 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2020, pp. 2981–2987.
- [9] K. M. Jatavallabhula, A. Kuwajerwala, Q. Gu, M. Omama, T. Chen, A. Maalouf, S. Li, G. Iyer, S. Saryazdi, N. Keetha, *et al.*, “Conceptfusion: Open-set multimodal 3d mapping,” *arXiv preprint arXiv:2302.07241*, 2023.
- [10] K. Weerakoon, A. J. Sathyamoorthy, J. Liang, T. Guan, U. Patel, and D. Manocha, “Graspe: Graph based multimodal fusion for robot navigation in outdoor environments,” *IEEE Robotics and Automation Letters*, 2023.
- [11] J. Wang, R. Lindenbergh, and M. Menenti, “Sigvox—a 3d feature matching algorithm for automatic street object recognition in mobile laser scanning point clouds,” *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 128, pp. 111–129, 2017.
- [12] H. Oleynikova, Z. Taylor, M. Fehr, R. Siegwart, and J. Nieto, “Voxblox: Incremental 3d euclidean signed distance fields for on-board mav planning,” in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2017, pp. 1366–1373.
- [13] S.-H. Zhang, Y.-C. Guo, and Q.-W. Gu, “Sketch2model: View-aware 3d modeling from single free-hand sketches,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 6012–6021.
- [14] B. Kerbl, G. Kopanas, T. Leimkühler, and G. Drettakis, “3d gaussian splatting for real-time radiance field rendering,” *ACM Trans. Graph.*, vol. 42, no. 4, pp. 139–1, 2023.
- [15] Y. Ji, Y. Liu, G. Xie, B. Ma, Z. Xie, and H. Liu, “Neds-slam: A neural explicit dense semantic slam framework using 3d gaussian splatting,” *IEEE Robotics and Automation Letters*, 2024.
- [16] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, “Nerf: Representing scenes as neural radiance fields for view synthesis,” in *ECCV*, 2020.
- [17] A. Yu, V. Ye, M. Tancik, and A. Kanazawa, “pixelnerf: Neural radiance fields from one or few images,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 4578–4587.
- [18] J. T. Barron, B. Mildenhall, M. Tancik, P. Hedman, R. Martin-Brualla, and P. P. Srinivasan, “Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 5855–5864.
- [19] M. Adamkiewicz, T. Chen, A. Caccavale, R. Gardner, P. Culbertson, J. Bohg, and M. Schwager, “Vision-only robot navigation in a neural radiance world,” *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 4606–4613, 2022.
- [20] J. J. Park, P. Florence, J. Straub, R. Newcombe, and S. Lovegrove, “Deepsdf: Learning continuous signed distance functions for shape representation,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 165–174.
- [21] Y.-C. Cheng, H.-Y. Lee, S. Tulyakov, A. G. Schwing, and L.-Y. Gui, “Sdfusion: Multimodal 3d shape completion, reconstruction, and generation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 4456–4465.
- [22] J. Engel, T. Schöps, and D. Cremers, “Lsd-slam: Large-scale direct monocular slam,” in *European conference on computer vision*. Springer, 2014, pp. 834–849.
- [23] R. Mur-Artal and J. D. Tardós, “Orb-slam2: An open-source slam system for monocular, stereo, and rgb-d cameras,” *IEEE transactions on robotics*, vol. 33, no. 5, pp. 1255–1262, 2017.
- [24] C. Campos, R. Elvira, J. J. G. Rodríguez, J. M. Montiel, and J. D. Tardós, “Orb-slam3: An accurate open-source library for visual, visual-inertial, and multimap slam,” *IEEE Transactions on Robotics*, vol. 37, no. 6, pp. 1874–1890, 2021.

- [25] 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 *2011 10th IEEE international symposium on mixed and augmented reality*. Ieee, 2011, pp. 127–136.
- [26] A. Dai, M. Nießner, M. Zollöfer, S. Izadi, and C. Theobalt, “Bundlefusion: Real-time globally consistent 3d reconstruction using on-the-fly surface re-integration,” *ACM Transactions on Graphics 2017 (TOG)*, 2017.
- [27] E. Sucar, S. Liu, J. Ortiz, and A. J. Davison, “imap: Implicit mapping and positioning in real-time,” in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 6229–6238.
- [28] Z. Zhu, S. Peng, V. Larsson, W. Xu, H. Bao, Z. Cui, M. R. Oswald, and M. Pollefeys, “Nice-slam: Neural implicit scalable encoding for slam,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 12 786–12 796.
- [29] H. Wang, J. Wang, and L. Agapito, “Co-slam: Joint coordinate and sparse parametric encodings for neural real-time slam,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 13 293–13 302.
- [30] Y.-L. Liu, C. Gao, A. Meuleman, H.-Y. Tseng, A. Saraf, C. Kim, Y.-Y. Chuang, J. Kopf, and J.-B. Huang, “Robust dynamic radiance fields,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 13–23.
- [31] G. S. Camps, R. Dyro, M. Pavone, and M. Schwager, “Learning deep sdf maps online for robot navigation and exploration,” *arXiv preprint arXiv:2207.10782*, 2022.
- [32] J. Sun, Y. Xie, L. Chen, X. Zhou, and H. Bao, “Neuralrecon: Real-time coherent 3d reconstruction from monocular video,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 15 598–15 607.
- [33] J. Ortiz, A. Clegg, J. Dong, E. Sucar, D. Novotny, M. Zollhoefer, and M. Mukadam, “isdf: Real-time neural signed distance fields for robot perception,” *arXiv preprint arXiv:2204.02296*, 2022.
- [34] B. Chen, R. Kwiatkowski, C. Vondrick, and H. Lipson, “Fully body visual self-modeling of robot morphologies,” *Science Robotics*, vol. 7, no. 68, p. eabn1944, 2022.
- [35] M. Wermelinger, P. Fankhauser, R. Diethelm, P. Krüsi, R. Siegwart, and M. Hutter, “Navigation planning for legged robots in challenging terrain,” in *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2016, pp. 1184–1189.
- [36] D. D. Fan, K. Otsu, Y. Kubo, A. Dixit, J. Burdick, and A.-A. Agha-Mohammadi, “Step: Stochastic traversability evaluation and planning for risk-aware off-road navigation,” *arXiv preprint arXiv:2103.02828*, 2021.
- [37] C. Cao, H. Zhu, F. Yang, Y. Xia, H. Choset, J. Oh, and J. Zhang, “Autonomous exploration development environment and the planning algorithms,” in *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022, pp. 8921–8928.
- [38] T. H. Y. Leung, D. Ignatyev, and A. Zolotas, “Hybrid terrain traversability analysis in off-road environments,” in *2022 8th International Conference on Automation, Robotics and Applications (ICARA)*. IEEE, 2022, pp. 50–56.
- [39] A. Shaban, X. Meng, J. Lee, B. Boots, and D. Fox, “Semantic terrain classification for off-road autonomous driving,” in *Proceedings of the 5th Conference on Robot Learning*, ser. Proceedings of Machine Learning Research, A. Faust, D. Hsu, and G. Neumann, Eds., vol. 164. PMLR, 08–11 Nov 2022, pp. 619–629. [Online]. Available: <https://proceedings.mlr.press/v164/shaban22a.html>
- [40] P. Roth, J. Nubert, F. Yang, M. Mittal, and M. Hutter, “Viplanner: Visual semantic imperative learning for local navigation,” *arXiv preprint arXiv:2310.00982*, 2023.
- [41] S. Jung, J. Lee, X. Meng, B. Boots, and A. Lambert, “V-strong: Visual self-supervised traversability learning for off-road navigation,” *arXiv preprint arXiv:2312.16016*, 2023.
- [42] X. Cai, M. Everett, J. Fink, and J. P. How, “Risk-aware off-road navigation via a learned speed distribution map,” in *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2022, pp. 2931–2937.
- [43] M. Mattamala, J. Frey, P. Libera, N. Chebrolu, G. Martius, C. Cadena, M. Hutter, and M. Fallon, “Wild visual navigation: Fast traversability learning via pre-trained models and online self-supervision,” *arXiv preprint arXiv:2404.07110*, 2024.
- [44] X. Cai, S. Ancha, L. Sharma, P. R. Osteen, B. Bucher, S. Phillips, J. Wang, M. Everett, N. Roy, and J. P. How, “Evora: Deep evidential traversability learning for risk-aware off-road autonomy,” *IEEE Transactions on Robotics*, 2024.
- [45] Z. Zhang, Q. Li, Z. Huang, J. Wu, J. Tenenbaum, and B. Freeman, “Shape and material from sound,” *Advances in Neural Information Processing Systems*, vol. 30, 2017.
- [46] B. Chen, M. Chiquier, H. Lipson, and C. Vondrick, “The boombox: Visual reconstruction from acoustic vibrations,” *arXiv preprint arXiv:2105.08052*, 2021.
- [47] J. Liu and B. Chen, “Sonicsense: Object perception from in-hand acoustic vibration,” in *8th Annual Conference on Robot Learning*, 2024. [Online]. Available: <https://openreview.net/forum?id=CpXiqz6qf4>
- [48] N. Riopelle, P. Caspers, and D. Sofge, “Terrain classification for autonomous vehicles using bat-inspired echolocation,” in *2018 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2018, pp. 1–6.
- [49] W. Mason, D. Brenken, F. Z. Dai, R. G. C. Castillo, O. S.-M. Cormier, and A. Sedal, “Acoustic tactile sensing for mobile robot wheels,” *arXiv preprint*

arXiv:2402.18682, 2024.

- [50] J. Christie and N. Kottege, “Acoustics based terrain classification for legged robots,” in *2016 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2016, pp. 3596–3603.
- [51] F. Xue, L. Hu, C. Yao, Z. Liu, Z. Zhu, and Z. Jia, “Sound-based terrain classification for multi-modal wheel-leg robots,” in *2022 international conference on advanced robotics and mechatronics (ICARM)*. IEEE, 2022, pp. 174–179.
- [52] G. B. Margolis and P. Agrawal, “Walk these ways: Tuning robot control for generalization with multiplicity of behavior,” in *Conference on Robot Learning*. PMLR, 2023, pp. 22–31.
- [53] M. Tancik, P. Srinivasan, B. Mildenhall, S. Fridovich-Keil, N. Raghavan, U. Singhal, R. Ramamoorthi, J. Barron, and R. Ng, “Fourier features let networks learn high frequency functions in low dimensional domains,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 7537–7547, 2020.
- [54] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, “Pointnet: Deep learning on point sets for 3d classification and segmentation,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 652–660.
- [55] V. Sitzmann, J. Martel, A. Bergman, D. Lindell, and G. Wetzstein, “Implicit neural representations with periodic activation functions,” *Advances in neural information processing systems*, vol. 33, pp. 7462–7473, 2020.