

# SALON: Self-supervised Adaptive Learning for Off-road Navigation

Matthew Sivaprakasam<sup>1</sup>, Samuel Triest<sup>1</sup>, Cherie Ho<sup>1</sup>, Shubhra Aich<sup>1</sup>,  
Jeric Lew<sup>2</sup>, Isaiah Adu<sup>3</sup>, Wenshan Wang<sup>1</sup>, and Sebastian Scherer<sup>1</sup>

**Abstract**—Autonomous robot navigation in off-road environments presents a number of challenges due to its lack of structure, making it difficult to handcraft robust heuristics for diverse scenarios. While learned methods using hand labels or self-supervised data improve generalizability, they often require a tremendous amount of data and can be vulnerable to domain shifts. To improve generalization in novel environments, recent works have incorporated adaptation and self-supervision to develop autonomous systems that can learn from their own experiences online. However, current works often rely on significant prior data, for example minutes of human teleoperation data for each terrain type, which is difficult to scale with more environments and robots. To address these limitations, we propose SALON, a perception-action framework for *fast* adaptation of traversability estimates with *minimal* human input. SALON rapidly learns online from experience while avoiding out of distribution terrains to produce adaptive and risk-aware cost and speed maps. Within *seconds* of collected experience, our results demonstrate comparable navigation performance over kilometer-scale courses in diverse off-road terrain as methods trained on 100-1000x more data. We additionally show promising results on significantly different robots in different environments. Our code is available at <https://theairlab.org/SALON>

## I. INTRODUCTION

Off-road autonomous driving is becoming an increasingly researched topic due to its wide range of applications. Robots are already being deployed in fields such as agriculture [1], infrastructure monitoring [2], and defense [3], where they must operate in unstructured and diverse environments. To perform reliably, they must reason about terrain lacking clear structures and markings to navigate from one goal to another without crashing or getting stuck.

Recent research has focused on improving navigation by generating costmaps with fine details, which captures the complexities of off-road terrain. For instance, detecting obstacles like rocks and trees without confusing them with traversable terrain such as small bushes and tall grass illustrates the level of detail required for effective navigation. Costmap generation methods, such as height-thresholding [4] and semantic segmentation [5,6] often struggle in off-road domain due to basic assumptions. More expressive geometric analyses [7–10] reduce this limitation, but can require extensive hand-tuning and still fail to distinguish different terrain with similar geometry. Learned methods have shown

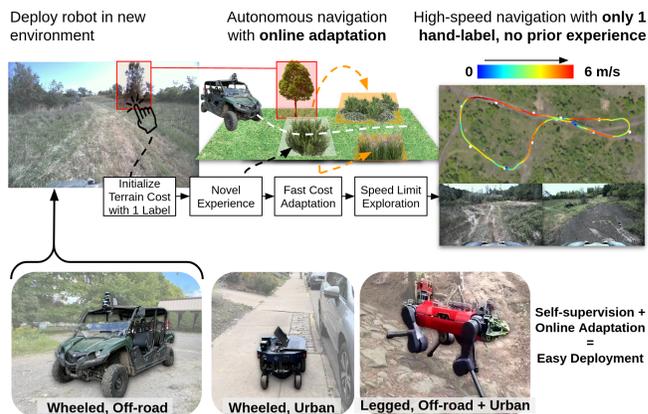


Fig. 1: We present SALON, a framework for off-road navigation with no prior experience. With one prior hand-label, our system running SALON learns from its own experience in the real-world to predict where and how fast to drive.

potential to address these issues [11–15] but are trained offline and struggle when deployed in new environments.

To improve generalization in novel environments, several works incorporate adaptation in a self-supervised manner, developing proprioceptive signals that enable the system to learn from its own experiences online [16–18]. *ALTER* [17] adapts its visual traversability model using explicitly generated cost (e.g. ideal to lethal) from LiDAR-built geometric maps, but is limited to accurate range sensors, and does not link to autonomy. Most relevant to our work, *WVN* [16] is a concurrent perception-action framework that learns traversability online by observing expert demonstration, then autonomously exploring within familiar terrain to refine its predictions. However, their anomaly detection severely penalizes areas where it has not traversed before, requiring teleoperation in all terrain types the robot should visit. This reliance on human teleoperation is not scalable when deploying in unknown environments, as the process may need to be repeated for different environments and robots. Such methods are insufficient for achieving *fast adaptation in new environments with minimal human input*.

We therefore propose SALON, a perception-action framework for fast adaptation of traversability estimates with minimal human input (only one click on an image, instead of minutes of expert teleoperation in all terrains). SALON achieves this with the following design choices:

- **Mapping and learning in map space:** We use a mapping pipeline to project visual features from the

\* This work was supported by ARL awards #W911NF1820218 and #W911NF20S0005.

<sup>1</sup> Robotics Institute, Carnegie Mellon University, [msivapra,striest,cherieh,saich,wenshanw,basti@andrew.cmu.edu](mailto:msivapra,striest,cherieh,saich,wenshanw,basti@andrew.cmu.edu)

<sup>2</sup> National University of Singapore, [jericlew@u.nus.edu](mailto:jericlew@u.nus.edu)

<sup>3</sup> Pennsylvania State University, [ioa5099@psu.edu](mailto:ioa5099@psu.edu)

camera into a Birds-Eye-View (BEV) map and associate the traversed cells with a cost and speed experienced by the robot. This results in cleaner, more distinctive maps.

- **One-shot cost augmentation:** a user simply clicks areas to avoid on prior images to initialize the system
- **Explicit Out-of-distribution (OOD) detection:** allows cleaner detection of anomalous objects, without long prior teleoperation in all safe terrain
- **Adaptive and risk-aware cost and speed maps:** Allows the robot to optimize mission-relevant metrics with novel experience

Overall we present three contributions:

- 1) SALON, a novel adaptive perception-action framework to generate **adaptive costmaps and speedmaps**, allowing autonomous navigation with minimal interventions, given as few as one hand-label.
- 2) **Real-world experiments** demonstrating performance at a similar level as state-of-the-art methods trained offline on 100-1000x more data.
- 3) Qualitative results on **multiple heterogeneous robots** differing in terms of dynamics, sensors, cost functions, visual back-ends, and environments.

Our code is open source, leveraging an existing mapping framework in order to facilitate deployment on other robots.

## II. RELATED WORK

### A. Costmap Generation for Off-road Driving

There exists a large number of existing works on learning self-supervised costmaps, both on and off-road. Some approaches leverage privileged information to supervise neural networks that predict map information at a given timestep [13,14,17,19], but this information consists of semantic segmentation or comes from handcrafted cost functions, both of which require copious amount of hand labels and tuning. Some methods aim to circumvent this explicit hand-labeling requirement by using expert demonstration data [11,20], and while the supervision comes from the data collection process itself new challenges arise with ensuring demonstration quality and adapting to novel stimuli without a human in the loop. Recent works, taking inspiration from older methods [7,21,22], have explored the potential for proprioception as supervision as it allows for a strong robot-specific relationship between experience and cost. Some of these methods leverage signals such as residuals between planned and expected trajectories [23,24] or IMU-based roughness score [12,25,26]. However, above methods require a significant amount of training data. While recent works have leveraged pre-trained models or visual foundation models (VFMs) to reduce the amount of labeled inputs [15,27], costmap generation methods that rely on statically-trained models may fail when deployed in out-of-distribution environments.

### B. Adaptive Methods for Costmap Generation

To improve performance when deployed in new environments, many works incorporate adaptation in their autonomy stack to learn from online experience. However, previous

works do not sufficiently achieve *fast adaptation in new environments with minimal human input* due to insufficiently expressive features, sensing restrictions, and need for minutes of pretraining data before effective adaptation. We find works that adapt online with LiDAR geometry measurements [17,28–31], but still struggle to reason about complex details present in natural environments. Instead, we propose a vision-based framework that is flexible to learn from different cost functions based on various sensors. We also find older works that learn from stereo camera measurements [32,33], but their use of less expressive image features hinders the prediction performance. More recently, adaptive methods have adopted the use of more expressive deep-learning based image features [16,18] to learn odometry-based cost, similar to our work. However, TerraPN [18] takes  $\sim 25$  minutes to learn which does not fulfill our need for fast adaptation, whereas ours can learn in the order of seconds of traversing a new type of terrain. Most relevant to our work, WVN [16] had success taking it a step further using Vision Transformer-based foundation model features to achieve faster, more generalizable adaptation. However, their framework requires human teleoperation in the different types of terrains the robot is expected to visit. Such need for human teleoperation is not scalable, as human teleoperation may be needed for each combination of environments and robots. In contrast, our method can deploy and explore new types of terrain with minimal human input (only one click on an image, instead of minutes of human teleoperation in all expected terrains).

## III. ADAPTING PERCEPTION WITH ONLINE EXPERIENCE

Our approach is powered by generalizable features, self-supervised cost signals, intelligent data management, and probabilistic traversability estimation (Fig. 2). Together, they enable our system to adapt quickly to novel terrain without extensive prior demonstrations.

### A. Visual Mapping

We leverage expressive features from VFMs aggregated into a BEV map as our terrain representation which is in turn used for prediction.

1) *Visual Backbone:* Pixel-level features in the image frame can be computed via a visual feature extractor, providing the mapping:

$$f_{\theta}(I_{3 \times H \times W}) = D \in \mathbb{R}^{C \times H \times W} \quad (1)$$

where  $I$  is an RGB image,  $D$  is a “featurized” image, with  $C$  feature channels generated by feature extractor  $f_{\theta}$ .

We use VFMs as feature extractors as we observe qualitatively that they not only can reason about commonly encountered terrain but also previously unexperienced terrain and objects. While the spatial resolution of these features is often lower due to model architecture, we find their powerful generalizability compensates for this shortcoming.

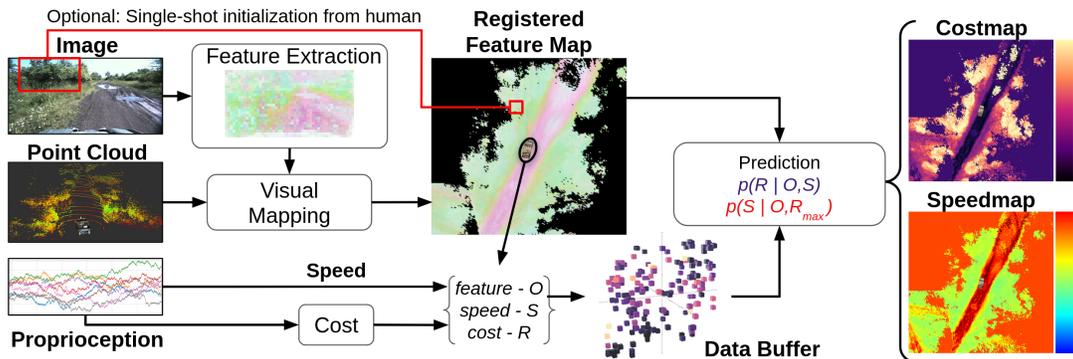


Fig. 2: SALON Overview: We learn rapidly from online experience with minimal human input to predict cost maps and speed maps. Visual foundation model (VFM) features in the map space, proprioceptive supervision, and smart data strategies together enable perception that adapts quickly to its environment.

2) *Dimensionality Reduction*: VFMs produce features with hundreds of channels per pixel, which can be intractable to represent in BEV in real-time. To compress them into a lower-dimensional space, we first generate feature images  $D$  from a subset of training images, then randomly sample pixel-level embeddings  $d \in \mathbb{R}^C$  from each image. Inspired by vector of locally aggregated descriptors (VLAD), which is popular in the place-recognition domain [34], we perform K-Means clustering on the training subset to generate  $k$  feature clusters ( $F_{1:k}$ ). A  $k$ -dimensional descriptor is then generated for each VFM pixel  $d$ , where the  $k$ -th feature in the descriptor is the  $L_1$  distance of  $d$  to cluster center  $F_k$  (Equation 2).

$$d_{VLAD}[k] = \|d - F_k\|_1 \quad (2)$$

In line with prior work [35], we find that this dimensionality reduction technique preserves semantic properties, which can provide strong priors for downstream tasks.

3) *BEV Mapping*: We aggregate these reduced FPV features in a BEV map using the same mapping method used in Velociraptor [15]. Using known calibration, visual features from a camera are associated with 3D points from a lidar, then projected into a BEV map which is aggregated over time by applying an exponential moving average.

### B. Curation of the Self-Supervised Signal

For a robot to adapt from its own experience, it requires a signal that associates different types of terrains with different costs in a way that matches human intuition. Following work by Castro et al. [12], we take inspiration from prior approaches [22,36,37] and use bandpower  $BP$  of a signal across a frequency range  $[f_{min}, f_{max}]$  as a way to compute roughness. While many works [12,28] compute bandpower for  $Z$ -axis (vertical) acceleration data alone, we find that for our full-scale system it is beneficial to include the other axes as well as readings from our vehicle suspension (henceforth referred to as shock travel).

In order to design a mapping from these sensor measurements to cost, we first collect a small dataset driving over different types of terrain at different speeds, periodically annotated by a passenger in the vehicle with a traversability

score in the range 0-1. The roughness is generated from computing the bandpower for all proprioceptive signals  $R = \sum_{i \in [a_x, a_y, a_z, shock \dots]} w_i BP(s_i, f_i^{min}, f_i^{max})$ , where  $w$  is the weight for each signal, and  $s$  is the window length of data. In order to obtain the best set of parameters  $[w_i, s_i, f_i^{min}, f_i^{max}]$  we optimize them to produce a roughness that matches the human annotations based on cumulative L1 error.

### C. Intelligent Data Maintenance

Leveraging the visual BEV mapping described above as a perceptual representation, we have the robot store experience as it drives, where a sample collected at time  $t$  contains:

- 1)  $O_t$  - The observed visual feature from the BEV map under the vehicle tire at time  $t$
- 2)  $S_t$  - The speed that the vehicle was traveling
- 3)  $R_t$  - The roughness that the vehicle experienced

Over time, the robot must throw out old samples to make room for new experiences. Rather than adopting a "first in, first out" (FIFO) strategy that can lead to catastrophic forgetting, we implement a strategy that aims to ensure an even distribution of data across the feature space. The VLAD features in the BEV representation by nature describe distance of observations to pre-defined clusters. Rather than throwing out the oldest sample, we leverage the semantic nature of our feature representation to instead throw out a sample  $n_{\bar{C}\bar{S}}$  corresponding to the most common "semantic class"  $\bar{C}$  and speed  $\bar{S}$ .

To verify this strategy, we compute the average pairwise distance between all points in the buffer for multiple sample trajectories, shown in Table I where the higher values for our method indicate a better coverage of the sample space in all of three different scenarios.

TABLE I: Avg. Pairwise Distance Between Points in Buffer

Strategy	Scenario 1	Scenario 2	Scenario 3
FIFO	3.47	3.60	3.59
Remove $n_{\bar{C}\bar{S}}$	<b>4.79</b>	<b>4.43</b>	<b>4.33</b>

#### D. Costmap and Speedmap Estimation

Given prior experience pairing roughness with visual features and speeds, the system needs a means of reasoning about the cost of the terrain ahead of it. We can predict the mean roughness  $\mu_R$  and variance  $v_R$  of a cell in the BEV map given its feature and the speed of the vehicle.

$$\mu_R, v_R = p(R|O, S) \quad (3)$$

The mean and variance can be computed using Gaussian Process Regression (GPR) using a radial basis function (RBF) kernel. Note that this can also be approximated with a simple MLP network, but we find that in our case GPR allows faster adaptation with stable predictions.

Note that this formulation also provides an estimate of the variance, which in turn means that the final roughness prediction can be tuned using Conditional Value at Risk (CVaR) based on the user’s preferred risk tolerance similar to other works [8,11,23]. The risk-adjusted predicted roughness assuming a Gaussian distribution becomes:

$$R = \mu_R + v_R \frac{\phi(\Phi^{-1}(\alpha_R))}{1 - \alpha_R} \quad (4)$$

Where  $\alpha_R$  is set by the user to vary risk-tolerance.

The experience buffer can also be similarly leveraged to predict speedmaps that dictate the upper-bound speed that the robot should travel for each cell. The user can specify a maximum desired roughness  $R_{max}$ , and use the same method to instead predict the mean speed  $\mu_S$  and variance  $v_S$ :

$$\mu_S, v_S = p(S|O, R = R_{max}), 0 \leq R_{max} \leq 1 \quad (5)$$

which can in turn be used to compute a speed limit for the downstream controller.

An issue arises here with out-of-distribution situations, where we are unlikely to obtain high-speed predictions without first experiencing high-speed data. We account for this by using CVaR with parameter  $\alpha_S$  similarly to  $\alpha_R$ , but dynamically adapting it instead of having it set by the user. While the vehicle is traveling within some margin of the speed limit but experiencing roughness that is significantly less than the expected roughness  $R_{max}$ ,  $\alpha_S$  is incrementally increased, and decreased if it is exceeding  $R_{max}$ . This allows the robot to explore higher speeds until it obtains evidence that supports its predictions, and adjust its limits if encounters previously unseen terrain that is too rough.

#### E. One-Shot Costmap Augmentation

Above, we show how a system can leverage features from VFMs to quickly learn costmaps and speedmaps without any human labels and adapt them with online experience. However, this formulation using signals such as roughness comes with the limitation that the robot can only learn about what it can physically drive on. This means reasonable predictions cannot be guaranteed for features that correspond to lethal objects.

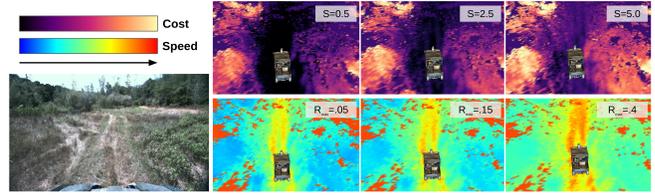


Fig. 3: Our system understands the relationship between velocity and traversability, and will command speeds that match the user’s tolerance for experienced cost. *Top row*: costmaps, conditioned on different speeds increasing from left to right; *Bottom row*: speedmaps, with a user-set maximum cost threshold increasing from left to right.

We address this problem by relaxing the assertion of zero human labels to one human label for commonly-experienced lethal objects, such as trees. We find that by simply choosing a single feature from the BEV map that corresponds to a tree and permanently associating it with high roughness in the buffer, we can obtain high cost values for all the trees that the robot experiences, without needing to train the network further or have a user label several trees.

1) *Avoiding Out-of-Distribution Terrain*: While the one-shot costmap augmentation is effective for frequently-encountered lethal terrain like trees, it may not be possible for a user to label all unique types of lethal terrain. In order to avoid terrain such as foreign objects, we leverage the same uncertainty estimation presented in Velociraptor, where a map cell is considered lethal if its distance to all feature clusters exceeds a threshold. We apply an additional morphological erosion and dilation to the uncertainty estimate in the BEV space to remove noise.

## IV. EXPERIMENTS AND RESULTS

### A. Platforms

To demonstrate the generalizability of our method, we show our results on three different robot platforms, with our perception running in-the-loop on our primary system and on teleoperated data from the other two. An overview is provided in Table II. For each robot we test a different combination of visual back-end and cost function.

1) *Full-Scale All-Terrain Vehicle*: We run our autonomy experiments on a Yamaha Viking All-Terrain Vehicle (ATV), first modified by Mai et al. [38] and further modified by Sivaprakasam et al. [39]. We use the ViT-B version of DINOv2 and use features from a held-out area to compute clusters for the VLAD descriptors.

2) *Autonomy-Equipped Wheelchair*: We test our method on data from a wheelchair equipped with a stereo camera driven in an urban environment, and use the ViT-B version of RADIOv2.5 [40]. We use the same roughness cost function

TABLE II: Robot Autonomy Configurations

Robot	Dynamics	Environment	Depth Sensor	Vis. Backbone	Cost Function
ATV	Ackermann	Nature	Lidar	DINOv2	Roughness
Wheelchair	Skid-Steer	Urban	Stereo Cam.	RADIOv2.5	Roughness
ANYmal	Quadruped	Urb. + Nat.	Lidar	DINOv2	Velocity Error

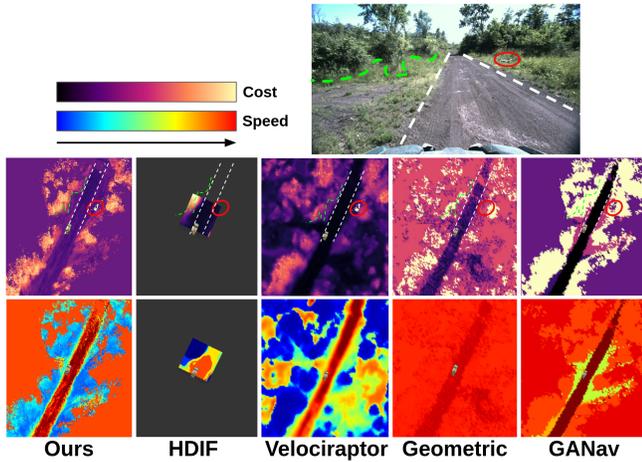


Fig. 4: Comparison of our method against baselines. Note our method’s ability to distinguish the tree line (green dashed line), trail (white dashed line), and the shattered TV hidden in the bushes (red circle).

used on the ATV (without the shock-travel sensor information), and one-shot cost augmentation is used to assign cars with high cost.

3) *ANYmal Quadruped*: We also test our method on an ANYmal quadruped robot developed by ANYbotics [41], using data provided from Wild Visual Navigation (WVN) [16]. We use the cost function provided in WVN, based on the discrepancy between desired and commanded velocity. One-shot cost augmentation to assigns trees and building walls with high cost. We compute descriptors using the same clusters generated from the off-road site used for the ATV.

### B. ATV Autonomous Navigation

1) *Baselines*: Using an MPPI [42] controller to actuate on the costmaps, we compare our method on our primary system against four other baselines:

- **Geometric cost function**: we use the cost function used in training ALTER [17] as a geometry-informed baseline. We compute a speedmap that encourages higher speed in areas with high planarity.
- **Semantic segmentation using GANav [43]**: we use GANav, mapping the semantic logits using the same

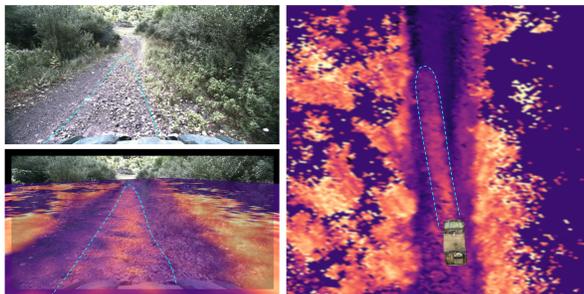


Fig. 5: SALON’s ability to distinguish fine-grained terrain: Rough gravel in the middle of the trail is high cost.

TABLE III: Capability Comparison Against Baselines

Method	Self Supervised	Online Adaptation	Risk Aware	OOD Detection	Velocity Informed
Geometry	✓	✗	✗	✗	✗
GA-Nav	✗	✗	✗	✗	✗
HDIF	✓	✗	✗	✗	✓
Velociraptor	✓	✗	✓	✓	✓
SALON	✓	✓	✓	✓	✓

pipeline used for our method, and associate the classes to hand-tuned costs and speeds.

- **How Does It Feel? (HDIF) [12]**: we compare against HDIF as it is a learning-based predecessor to this work, using a constant desired speed and limited range due to runtime constraints.
- **Velociraptor [15]**: an inverse reinforcement learning approach that leverages geometric and visual cues to predict costmaps and speedmaps.

2) *Qualitative Results*: Within less than 60 seconds of experience, our method appears to outperform the simpler baselines (GANav and geometric cost function) and predict maps at a level of detail more similar to Velociraptor (Fig. 4). Note that Velociraptor also uses geometric features from lidar as an input, enabling more range and a wider field-of-view.

We find that after adaptation, our predicted maps reflect expected behavior conditioned on various velocities and max roughness values (Fig. 3). We also present additional qualitative examples that highlight our method’s ability to not only distinguish high-level terrain types — such as trees, grass, and trail — but also pick up on visual cues such as grass density and gravel in order to make expressive predictions (Fig. 5).

To better observe the adaptive behavior of our method, we show predictions less than 10 seconds apart in the same lap (Fig. 6). Without any experience on grass, our method can pick out trees and trail, but believes that the grass is a higher cost than it should be. After driving on it for just a few seconds it has a finer understanding of the terrain, able to distinguish between smooth and rough vegetation.



Fig. 6: Example of SALON’s fast adaptation: Within 10 seconds of experiencing grass, SALON is able to quickly differentiate key terrains, such as, ideal short grass, riskier vegetation and lethal trees.

TABLE IV: Navigation Performance Metrics

Method	# Interactions					# Undesirable Behavior					Avg. Speed (m/s)					Avg. Roughness				
	Lap Number	1	2	3	4	5	1	2	3	4	5	1	2	3	4	5	1	2	3	4
Geometry	3	6	5	5	4	1	1	2	1	2	<b>5.2</b>	4.5	4.5	4.6	4.6	.30	.26	.27	.27	.27
GANav	4	6	4	4	4	1	2	2	2	2	3.7	3.5	3.5	3.7	3.6	.18	.19	.18	<b>.19</b>	.19
HDIF	3	6	3	3	4	<b>0</b>	1	2	1	<b>0</b>	4.2	3.6	4.0	4.2	3.8	.23	.21	.24	.23	.23
Velociraptor	1	<b>0</b>	1	1	1	1	2	2	1	2	5.1	<b>4.7</b>	4.7	<b>4.7</b>	<b>4.9</b>	.26	.26	.28	.27	.28
SALON, $R_{max} = .2$	2	1	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	1	<b>0</b>	1	<b>0</b>	3.4	3.6	3.4	3.7	3.9	<b>.16</b>	<b>.17</b>	<b>.17</b>	<b>.19</b>	<b>.18</b>
SALON, $R_{max} = .4$	2	<b>0</b>	1	<b>0</b>	1	<b>0</b>	<b>0</b>	1	1	<b>0</b>	3.8	4.5	<b>4.8</b>	4.6	4.5	.19	.25	.25	.24	.24

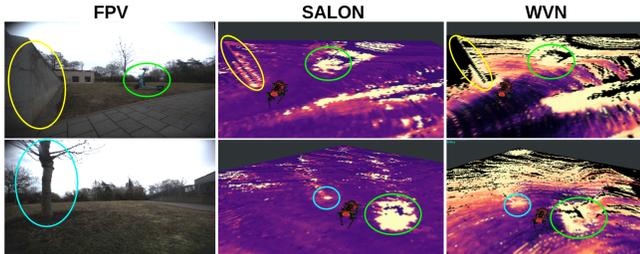


Fig. 7: With the same amount of data as WVN, our method is able to correctly cost lethal objects like trees and walls without incorrectly costing short grass.

3) *Quantitative Results:* We run all baselines as well as our own method for five laps in a “figure 8” style course with 50m waypoint spacing, evaluated on four metrics (Table IV). We count the number of times the safety driver intervened in order to prevent damage, as well as the number of “undesirable behaviors” where the driver didn’t need to intervene but the system could have taken a better route (for example driving through a rough patch of grass when there is a trail right next to it).

Our method outperforms the baselines in all metrics apart from average speed. While Velociraptor had a higher average speed it also had more interventions, some of which were a result of driving too fast. Due to the fragility of our vehicle we prioritize number of interventions over speed, but we recognize that there is an inherent trade-off between the two and preferences may vary based on robot and operator.

### C. Qualitative Comparison to State-of-the-Art Adaptation Method

We test our method on the quadruped used in WVN. To compare the two methods qualitatively, we project the FPV cost predictions from WVN into a costmap using the same mapping pipeline in our approach. Our method appears to be more consistent as shown in Fig. 7, likely due to the fact that our method makes predictions in the map space, in contrast to than mapping predictions made in the FPV space. Additionally, while WVN appears to cost grass more heavily than pavement compared to ours we find that this is not representative of their cost function, which considers the two types of terrain to be similar.

### D. Qualitative Evaluation in Urban Environment

Using the same cost function used on the ATV but a different visual backend (RADIOv2.5), we run our method



Fig. 8: Evaluation on urban wheelchair: After driving over rough cobblestone, the system quickly recognizes within 5 seconds that it is much rougher than the smooth sidewalk.

on a wheelchair on sidewalks in an urban environment. In Fig. 8, we highlight a scenario where initially, the robot incorrectly equates rough cobblestone with smooth pavement. After driving over the rough section for a few seconds, it learns to distinguish between the two in the environment ahead of it. We also find that it is able to pick up on fine details such as cracks in the sidewalk.

## V. CONCLUSIONS AND FUTURE WORK

In this work we present SALON, a framework for predicting costmaps and speedmaps that allow a system to adapt in real time to novel experiences by associating generalizable features from visual foundation models with proprioceptive signals. This results in a system that can make more nuanced predictions about its environment than the prior state-of-the-art that in turn allow for improved navigation behaviors. Further, we demonstrate promising results on multiple robots to highlight the generalizability of our method.

While we improve on many of the issues presented in our baselines, there still remains a number of future directions. For example, the exploration of the trade-offs between different distribution assumptions would be highly beneficial. Additionally, our strategy for speedmap prediction assumes that roughness increases relatively monotonically with speed. While we observe this to be mostly the case in our system, there could exist other systems built in a way such that certain terrain is actually less rough at high speeds.

There is also a need for unified benchmarking of off-road autonomy as a whole. While signals such as interventions and speed are certainly strong identifiers of success, they often have an inverse relationship with the weight of each being ambiguous due to vehicle and user constraints.

## ACKNOWLEDGEMENT

We thank Humeyra Kacar and Anton Yanovich for their assistance in field testing. We'd also like to thank Nathan Litzinger, Bangjie Xue, Victor Zayakov, and Jiahe Xu for their work on the setting up the wheelchair testing platform for data collection.

## REFERENCES

- [1] D. F. Yépez-Ponce, J. V. Salcedo, P. D. Rosero-Montalvo, and J. Sanchis, "Mobile robotics in smart farming: current trends and applications," *Frontiers Artif. Intell.*, vol. 6, 2023. [Online]. Available: <https://doi.org/10.3389/fraci.2023.1213330>
- [2] S. Gibb, H. M. La, T. Le, L. Nguyen, R. Schmid, and H. Pham, "Nondestructive evaluation sensor fusion with autonomous robotic system for civil infrastructure inspection," *Journal of Field Robotics*, vol. 35, no. 6, pp. 988–1004, 2018. [Online]. Available: <https://onlinelibrary.wiley.com/doi/abs/10.1002/rob.21791>
- [3] J. E. Naranjo, M. Clavijo, F. Jiménez, O. Gómez, J. L. Rivera, and M. Anguita, "Autonomous vehicle for surveillance missions in off-road environment," in *2016 IEEE Intelligent Vehicles Symposium (IV)*, 2016, pp. 98–103.
- [4] A. Elfes, "Using occupancy grids for mobile robot perception and navigation," *Computer*, vol. 22, pp. 46 – 57, 07 1989.
- [5] D. Maturana, P.-W. Chou, M. Uenoyama, and S. Scherer, "Real-time semantic mapping for autonomous off-road navigation," in *Field and Service Robotics*, M. Hutter and R. Siegwart, Eds. Cham: Springer International Publishing, 2018, pp. 335–350.
- [6] 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>
- [7] P. Fankhauser, M. Bjelonic, C. Dario Bellicoso, T. Miki, and M. Hutter, "Robust rough-terrain locomotion with a quadrupedal robot," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*, 2018, pp. 5761–5768.
- [8] A. Dixit, D. Fan, K. Otsu, S. Dey, A.-A. Agha-Mohammadi, and J. Burdick, "Step: Stochastic traversability evaluation and planning for risk-aware off-road navigation; results from the darpa subterranean challenge," *Field Robotics*, vol. 4, no. 1, p. 182–210, Jan. 2024. [Online]. Available: <http://dx.doi.org/10.55417/fr.2024006>
- [9] P. Krüsi, P. Furgale, M. Bosse, and R. Siegwart, "Driving on point clouds: Motion planning, trajectory optimization, and terrain assessment in generic nonplanar environments: Driving on point clouds," *Journal of Field Robotics*, vol. 34, 12 2016.
- [10] C. Noren, B. Vundurthy, S. Scherer, and M. Travers, "Interaction-aware control for robotic vegetation override in off-road environments," *Journal of Terramechanics*, vol. 117, p. 101034, 2025. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0022489824000764>
- [11] S. Triest, M. G. Castro, P. Maheshwari, M. Sivaprakasam, W. Wang, and S. Scherer, "Learning risk-aware costmaps via inverse reinforcement learning for off-road navigation," 2023.
- [12] M. G. Castro, S. Triest, W. Wang, J. M. Gregory, F. Sanchez, J. G. Rogers, and S. Scherer, "How does it feel? self-supervised costmap learning for off-road vehicle traversability," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*, 2023, pp. 931–938.
- [13] X. Meng, N. Hatch, A. Lambert, A. Li, N. Wagener, M. Schmittle, J. Lee, W. Yuan, Z. Chen, S. Deng, G. Okopal, D. Fox, B. Boots, and A. Shaban, "Terrainnet: Visual modeling of complex terrain for high-speed, off-road navigation," 2023.
- [14] J. Frey, S. Khattak, M. Patel, D. Atha, J. Nubert, C. Padgett, M. Hutter, and P. Spieler, "Roadrunner - learning traversability estimation for autonomous off-road driving," 2024. [Online]. Available: <https://arxiv.org/abs/2402.19341>
- [15] S. Triest, M. Sivaprakasam, S. Aich, D. Fan, W. Wang, and S. Scherer, "Velociraptor: Leveraging visual foundation models for label-free, risk-aware off-road navigation," in *8th Annual Conference on Robot Learning*, 2024. [Online]. Available: <https://openreview.net/forum?id=AhEE5wrcLU>
- [16] M. Mattamala, J. Frey, P. Libera, N. Chebrolo, G. Martius, C. Cadena, M. Hutter, and M. Fallon, "Wild visual navigation: Fast traversability learning via pre-trained models and online self-supervision," 2024. [Online]. Available: <https://arxiv.org/abs/2404.07110>
- [17] E. Chen, C. Ho, M. Maulimov, C. Wang, and S. Scherer, "Learning-on-the-drive: Self-supervised adaptation of visual offroad traversability models," 2023.
- [18] A. J. Sathyamoorthy, K. Weerakoon, T. Guan, J. Liang, and D. Manocha, "Terrappn: Unstructured terrain navigation using online self-supervised learning," in *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2022, pp. 7197–7204.
- [19] S. Triest, D. D. Fan, S. Scherer, and A.-A. Agha-Mohammadi, "Unrealnet: Learning uncertainty-aware navigation features from high-fidelity scans of real environments," 2024. [Online]. Available: <https://arxiv.org/abs/2407.08720>
- [20] M. Wulfmeier, D. Rao, D. Z. Wang, P. Ondruska, and I. Posner, "Large-scale cost function learning for path planning using deep inverse reinforcement learning," *The International Journal of Robotics Research*, vol. 36, no. 10, pp. 1073–1087, 2017. [Online]. Available: <https://doi.org/10.1177/0278364917722396>
- [21] A. Angelova, L. H. Matthies, D. M. Helmick, and P. Perona, "Learning and prediction of slip from visual information," *Journal of Field Robotics*, vol. 24, 2007. [Online]. Available: <https://api.semanticscholar.org/CorpusID:15679808>
- [22] E. DuPont, C. Moore, E. Collins, and E. Coyle, "Frequency response method for terrain classification in autonomous ground vehicles," *Autonomous Robots*, vol. 24, pp. 337–347, 05 2008.
- [23] 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," 2024. [Online]. Available: <https://arxiv.org/abs/2311.06234>
- [24] M. Sivaprakasam, S. Triest, W. Wang, P. Yin, and S. Scherer, "Improving off-road planning techniques with learned costs from physical interactions," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*, 2021, pp. 4844–4850.
- [25] A. Pokhrel, A. Datar, M. Nazeri, and X. Xiao, "Cahsor: Competence-aware high-speed off-road ground navigation in se(3)," 2024. [Online]. Available: <https://arxiv.org/abs/2402.07065>
- [26] X. Yao, J. Zhang, and J. Oh, "Rca: Ride comfort-aware visual navigation via self-supervised learning," in *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2022, pp. 7847–7852.
- [27] S. Jung, J. Lee, X. Meng, B. Boots, and A. Lambert, "V-strong: Visual self-supervised traversability learning for off-road navigation," 2024. [Online]. Available: <https://arxiv.org/abs/2312.16016>
- [28] J. Seo, T. Kim, S. Ahn, and K. Kwak, "Metaverse: Meta-learning traversability cost map for off-road navigation," 2024. [Online]. Available: <https://arxiv.org/abs/2307.13991>
- [29] H. Dahlkamp, A. Kaehler, D. Stavens, S. Thrun, and G. R. Bradski, "Self-supervised monocular road detection in desert terrain." in *Robotics: science and systems*, vol. 38. Philadelphia, 2006.
- [30] J. A. Bagnell, D. Bradley, D. Silver, B. Sofman, and A. Stentz, "Learning for autonomous navigation," *IEEE Robotics and Automation Magazine*, vol. 17, no. 2, pp. 74–84, 2010.
- [31] B. Sofman, E. Lin, J. A. Bagnell, N. Vandapel, and A. Stentz, "Improving robot navigation through self-supervised online learning," in *Robotics: Science and Systems*, 2006.
- [32] R. Hadsell, A. Erkan, P. Sermanet, M. Scoffier, U. Muller, and Y. LeCun, "Deep belief net learning in a long-range vision system for autonomous off-road driving," in *2008 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2008, pp. 628–633.
- [33] M. Happold, M. Ollis, and N. Johnson, "Enhancing supervised terrain classification with predictive unsupervised learning," in *Proceedings of Robotics: Science and Systems*, Philadelphia, USA, August 2006.
- [34] H. Jégou, M. Douze, C. Schmid, and P. Pérez, "Aggregating local descriptors into a compact image representation," in *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2010, pp. 3304–3311.
- [35] N. Keetha, A. Mishra, J. Karhade, K. M. Jatavallabhula, S. Scherer, M. Krishna, and S. Garg, "Anyloc: Towards universal visual place recognition," *IEEE Robotics and Automation Letters*, vol. 9, no. 2, pp. 1286–1293, 2023.
- [36] D. Stavens and S. Thrun, "A self-supervised terrain roughness

- estimator for off-road autonomous driving,” 2012. [Online]. Available: <https://arxiv.org/abs/1206.6872>
- [37] L. Wellhausen, A. Dosovitskiy, R. Ranftl, K. Walas, C. Cadena, and M. Hutter, “Where should i walk? predicting terrain properties from images via self-supervised learning,” *IEEE Robotics and Automation Letters*, vol. 4, no. 2, pp. 1509–1516, 2019.
- [38] J. Mai, “System design, modelling, and control for an off-road autonomous ground vehicle,” Master’s thesis, Carnegie Mellon University, Pittsburgh, PA, July 2020.
- [39] M. Sivaprakasam, P. Maheshwari, M. G. Castro, S. Triest, M. Nye, S. Willits, A. Saba, W. Wang, and S. Scherer, “Tartandrive 2.0: More modalities and better infrastructure to further self-supervised learning research in off-road driving tasks,” in *2024 IEEE International Conference on Robotics and Automation (ICRA)*, 2024, pp. 12 606–12 606.
- [40] M. Ranzinger, G. Heinrich, J. Kautz, and P. Molchanov, “Am-radio: Agglomerative vision foundation model reduce all domains into one,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2024, pp. 12 490–12 500.
- [41] M. Hutter, C. Gehring, D. Jud, A. Lauber, C. D. Bellicoso, V. Tsounis, J. Hwangbo, K. Bodie, P. Fankhauser, M. Bloesch, R. Diethelm, S. Bachmann, A. Melzer, and M. Hoepflinger, “Anymal - a highly mobile and dynamic quadrupedal robot,” in *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2016, pp. 38–44.
- [42] G. Williams, A. Aldrich, and E. A. Theodorou, “Model predictive path integral control: From theory to parallel computation,” *Journal of Guidance, Control, and Dynamics*, vol. 40, no. 2, pp. 344–357, 2017. [Online]. Available: <https://doi.org/10.2514/1.G001921>
- [43] T. Guan, D. Kothandaraman, R. Chandra, A. J. Sathyamoorthy, K. Weerakoon, and D. Manocha, “Ga-nav: Efficient terrain segmentation for robot navigation in unstructured outdoor environments,” *IEEE Robotics and Automation Letters*, vol. 7, no. 3, pp. 8138–8145, 2022.