

VANP: Learning Where to See for Navigation with Self-Supervised Vision-Action Pre-Training

Mohammad Nazeri, Junzhe Wang, Amirreza Payandeh, and Xuesu Xiao

Abstract—Humans excel at efficiently navigating through crowds without collision by focusing on specific visual regions relevant to navigation. However, most robotic visual navigation methods rely on deep learning models pre-trained on vision tasks, which prioritize salient objects—not necessarily relevant to navigation and potentially misleading. Alternative approaches train specialized navigation models from scratch, requiring significant computation. On the other hand, self-supervised learning has revolutionized computer vision and natural language processing, but its application to robotic navigation remains underexplored due to the difficulty of defining effective self-supervision signals. Motivated by these observations, in this work, we propose a Self-Supervised Vision-Action Model for Visual Navigation Pre-Training (VANP). Instead of detecting salient objects that are beneficial for tasks such as classification or detection, VANP learns to focus only on specific visual regions that are relevant to the navigation task. To achieve this, VANP uses a history of visual observations, future actions, and a goal image for self-supervision, and embeds them using two small Transformer Encoders. Then, VANP maximizes the information between the embeddings by using a mutual information maximization objective function. We demonstrate that most VANP-extracted features match with human navigation intuition. VANP achieves comparable performance as models learned end-to-end with half the training time and models trained on a large-scale, fully supervised dataset, i.e., ImageNet, with only 0.08% data.

I. INTRODUCTION

In recent years, imitation learning, particularly behavior cloning [1], has become a leading approach for visual navigation models [2]–[8]. However, the performance of these models heavily relies on the visual features extracted by the model’s visual encoder. Although the limited memory and processing power onboard robots restrict the size of models deployable in real time, with such limitations we still need accurate and efficient onboard visual encoders, making convolutional neural networks (CNNs) more desirable than larger Vision Transformer models (ViTs) [9].

Training a visual navigation-specific encoder from scratch requires a large amount of data, leading to high computational demands and extended training times [10], [11]. To reduce this computational burden, most approaches use pre-trained vision models [4], [5]. While these models provide a decent scene representation, they specialize in extracting salient features for vision tasks such as object classification and detection [12]. These features may not always align with what is crucial for navigation [13]. For example, following sidewalks, avoiding grass, or navigating around stairs and

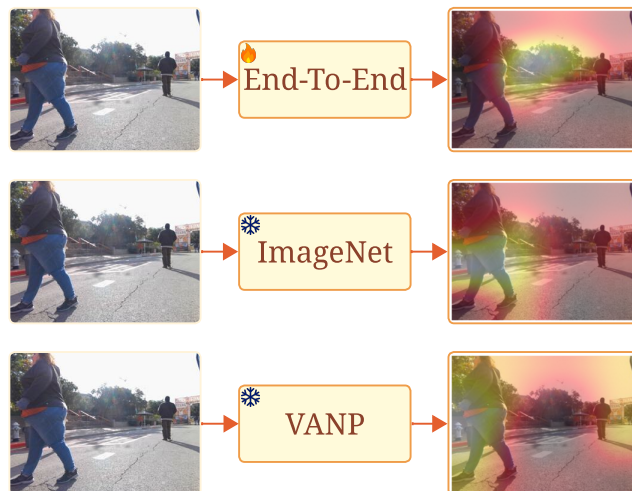


Fig. 1: Comparison of Activation Maps Learned by End-to-End, ImageNet, and VANP. VANP can extract multiple regions of interest for navigation without downstream navigation supervision compared to single salient regions by End-to-End and ImageNet pre-trained models.

guardrails are essential for robots, but these features might not be captured by encoders trained for generic vision tasks. Consequently, pre-trained models, like those trained on ImageNet, can sometimes lead to navigation failures by focusing on irrelevant distractions [10], [11].

Self-Supervised Learning (SSL) [14]–[17] has shown success in various computer vision tasks by extracting general features adaptable to downstream tasks with/without fine-tuning. For instance, Deep Neural Networks (DNNs) can be trained to predict the rotation of an image [18] or to reconstruct an image from its corrupted/obstructed version [19]. By completing these pretext tasks, DNNs learn to extract meaningful features from the data, which can be used to solve downstream tasks such as image classification and object detection [20]. However, a discrepancy exists between features extracted from generic models and those specifically needed for navigation. This leads us to ask the question: *can we train visual encoders that extract only navigation-relevant features using self-supervision?*

Considering both the success of SSL on a variety of computer vision tasks and the oftentimes mismatched features provided by generic SSL models for navigation tasks, we present Vision-Action Navigation Pretraining (VANP),

a non-contrastive self-supervised approach that completely relies on a navigation-specific pretext task to train the visual encoder without the need for negative samples.

The core idea behind VANP is inspired by how humans navigate in crowded spaces. We do not need to pay attention to all the people and objects in the scene, but only the ones that affect our navigation trajectory. To this end, VANP embeds visual history, future actions, and visual goal as self-supervision signals and leverages Transformers with additional context tokens (inspired by Bert [21] and VisionTransformers [9]) to generate embeddings. Then, VANP utilizes VICReg [23] as the pretext objective function to maximize the mutual information between the embeddings. The trained visual encoder can therefore discard redundant features unnecessary for navigation and focus only on navigation-relevant regions. For example, Fig. 1 shows the activation map of the last layer of ResNet-50 [24] trained with different methods. VANP learns navigation-relevant visual features with the help of our navigation-specific self-supervision signals.

Our experimental results suggest that VANP-extracted features trained on a dataset [25] that only contains 0.08% samples compared to ImageNet are as informative for a downstream navigation task as using ImageNet features. The contributions of this work can be summarized as follows:

- An SSL framework to train a visual encoder for robotic navigation tasks;
- Insights into what is happening inside CNNs during navigation using different approaches; and
- A benchmark on short and long-term navigation interaction to show the performance of different approaches.

II. RELATED WORK

Recent advances in natural language processing and computer vision, particularly those driven by self-supervised learning (SSL), motivate our work. In this section, we first compare common SSL approaches and then categorize applications of SSL into two groups for robotics and review their related works.

Self-Supervised Learning: SSL has shown promising results in recent years by almost reaching the performance of supervised baselines [15], [23]. Within SSL, two primary approaches have emerged: contrastive methods and information maximization methods. Both methodologies benefit from the use of the Siamese network architecture [26]. Contrastive methods [14], [27] typically require large data batches and leverage loss functions designed to explicitly push dissimilar data points away from the representations of similar data. Consequently, the performance of these methods is highly dependent on the quality and quantity of negative samples [28]. Recent advancements have led to the development of contrastive approaches that do not necessitate negative samples for learning effective embeddings. These methods employ various strategies to achieve comparable performance, such as BYOL [29], which utilizes a momentum encoder where one head receives a low-pass filtered version of the other. Alternatively, SimSiam [30] achieves

similar results by halting gradient flow within one of the heads.

Information maximization methods such as BarlowTwins [15] maximize the information between two heads by enforcing the empirical cross-correlation between the embeddings of both heads to be equivalent to the identity matrix. Additionally, VICReg [23] incorporates regularization terms to prevent information collapse, particularly in scenarios involving multimodal data. Therefore, VANP leverages VICReg to learn visual features by maximizing the information between different modalities.

Pre-training for Better Representation: Codevilla *et al.* [4] demonstrated the value of pre-trained models for training better policies in autonomous vehicles. Subsequently, many works adopted pre-trained computer vision models, often trained on ImageNet [4]–[7], [22], [31], [32]. However, general-purpose “foundation models” pre-trained on pretext tasks can achieve richer representations, enabling them to generalize to various downstream tasks with minimal data in a zero- or few-shot manner [20], [33], [34].

The literature has extensively studied foundation models for robot manipulation [35]–[40]. For example, R3M [41] pre-trained a general visual encoder for manipulation tasks on the Ego4D human video dataset [42], while CLIPort [43] leveraged the CLIP model [44] to enable language instructions for manipulation. Dadashi *et al.* proposed AQUA-Dem [35], a framework to learn quantized actions from demonstrations in continuous action spaces, while VANP is doing the opposite by learning visual features from continuous action spaces. Luo *et al.* [36] improved AQUA-Dem by using VQ-VAE [45] for offline reinforcement learning. Huang *et al.* [37] proposed Skill Transformer to learn long-horizon robotic tasks with the help of Transformers [46].

Inspired by Taskonomy [47], Shen *et al.* proposed conditioning visual demonstrations like segmentation and depth maps on actions during fusion rather than employing a naive fusion approach [48]. Yang *et al.* [49] projected the visual cues for navigation on the image space and then trained a policy on the augmented image. STERLING [50] and CAHSOR [51] have explored the concept of human preference learning and competence-awareness in the context of off-road navigation using SSL. These methods aligned sensor and visual embeddings by maximizing the mutual information between embeddings by leveraging VICReg [23] and BarlowTwins [15] respectively.

The work by Eftekhari *et al.* [52] presented the closest approach to VANP, employing a learnable codebook module to selectively filter visual observations based on the specific task. However, relying on task-relevant information, e.g., picking up the key, requires additional information that is not available without human annotation or using a simulator while VANP does not need access to such information to learn visual features. Wang *et al.* used noise to pre-train a visual encoder by predicting the scale of a patch within the noise image that applies to crop the goal observable from the current frame in real experiments [53]. In contrast, VANP deliberately disregards such task-specific information, focus-



Fig. 2: **VANP Architecture.** VANP learns to embed temporal features into spatial features by using a sequence of images and leveraging two TransformerEncoders with context tokens. VANP’s loss maximizes the mutual information between history, future actions, and the goal (left). Then, by appending an MLP to the Transformer context token, VANP predicts future trajectories during the downstream navigation task (right).

ing instead on extracting general navigation-relevant features. Another work closely related to VANP is NavFormer [54], which utilized BYOL [29] on two input images retrieved from a simulator. These images differ in the presence of dynamic objects within the scene. However, this approach confines NavFormer to the simulated environment, limiting its applicability to simulation environments where we have full control of the environment, e.g., making objects invisible to learn the importance of the presence and absence of the object as an obstacle. Conversely, VANP achieves real-world data generalization without relying on the pre-definition of specific rules only possible in simulation or through human annotation.

Pre-training for Better Policies: Foundation models hold promise for learning not only rich representations but also policies that can generalize across robotic tasks. For instance, SayCan [55] integrates pre-trained language skills with robot actions, demonstrating the potential of pre-training for robotic tasks. This allows robots to physically execute tasks, while the language model provides high-level task insights. Evaluations of real-world robotic tasks confirm the effectiveness of this grounded approach in handling abstract, long-duration instructions for a mobile manipulator. Li *et al.* [56] pre-trained language models to initialize policy networks predicting actions. Reid *et al.* [57] fine-tuned pre-trained sequence models on offline reinforcement learning tasks as the policy backbone. VPT [58] used pseudo-labeled Minecraft YouTube videos to learn a behavior cloning policy that can craft diamonds. VPT learns the inverse dynamics while VANP uses dynamics to learn visual features. GNM [59] learned a general policy to drive any robot by combining multiple datasets of different robot types. ViNT [60] further improved GNM by replacing the policy network with a Transformer [46].

III. METHODOLOGY

Learning visual features for robot navigation using only RGB camera input presents several challenges. Unlike traditional approaches that rely on LiDAR or depth cameras, RGB cameras lack explicit geometric information, making

navigation more complex [61]–[66]. Here, we formally define the visual navigation task and the learning setting for Vision-Action Navigation Pre-training (VANP).

A. Problem Definition

We define visual navigation as the task of navigating an environment with only RGB camera input, as explored in previous works [3], [4], [7]. The visual navigation problem can be formalized as follows. **Input:** The robot is given a sequence of past and current images from its front-facing camera, $o_t = [I_{t-\tau_P}, I_{t-\tau_P+1}, \dots, I_t] \in \mathcal{O}$, where t is the current time step, τ_P is the number of past frames, and \mathcal{O} is the space of all possible image sequences. The robot is also given its current goal e.g., GPS coordinates, pose, image, or next local coordinate in 2D space, $g \in \mathcal{G}$, which determines the direction it should move in the next time step. **Output:** The robot must select an action $a_t \in \mathcal{A}$ consisting of continuous linear and angular velocities. $\mathcal{A} = [-1, 1]^2$ is the action space, where $[-1, 1]$ maps to the minimal and maximal linear and angular velocity of the robot. **Visual Navigation:** The goal is to learn a policy, $\pi_\theta : \mathcal{O} \times \mathcal{G} \rightarrow \mathcal{A}$, where θ represents the policy’s parameters, to determine which action to take at each time step to reach its goal efficiently while avoiding collisions with others.

End-To-End models: For end-to-end or holistic models, we define the policy π_θ as follows: $a = \pi_\theta(o, g) = \sigma_\zeta(p_\phi(o) \oplus q_\psi(g))$, where σ is the controller policy parameterized by ζ , p is the image encoder parameterized by ϕ , q is the goal encoder parameterized by ψ , and \oplus is the aggregation of two vectors. To learn these parameters, two common approaches are (1) to learn all of them together in an end-to-end manner which makes the training difficult and time-consuming or (2) to pre-train the image encoder separately and only fine-tune the goal encoder along with the controller to reduce training time.

Challenges in visual feature learning: While extensive research has explored learning visual features for computer vision tasks using SSL [14], [15], [23], [27], [29], [30], adapting these models to specific tasks presents unique challenges [50], [51]. Images in the real world contain implicit cues for navigation but are sometimes full of re-

dundant information. In the context of visual navigation, one such challenge lies in learning visual features from image sequences without unnecessarily capturing such a redundancy, which may result in ambiguity. Additionally, it is not trivial to extract contrastive learning signals from visual navigation actions for contrastive SSL, e.g., an action appropriate for one scenario may or may not be appropriate for another, or different actions may be appropriate for the same scenario. For instance, in a scenario where a pedestrian stands in front of the robot, two equally valid actions exist: overtaking from either the left or right side. In such cases, simply negating the angular velocity cannot yield a meaningful negative sample and can introduce ambiguity. Furthermore, employing actions from different sequences as negative samples might not provide pertinent information for visual navigation, as actions are inherently influenced by the observed environment. In the next section, we show how VANP addresses these challenges and trains the image encoder p without a downstream objective function.

B. Vision-Action Model

VANP leverages VICReg [23] to maximize the information between past observations, a future goal, and future actions while maintaining the information collapse between input heads to train the image encoder p . Unlike vision SSL models that work on the joint embedding of augmented images [27], [67], VANP correlates the action space \mathcal{A} and goal space \mathcal{G} with the pixel latent space \mathcal{O} as shown in Fig. 2. We define VANP pre-training as follows: We sample a batch of $(I_{t-\tau_P:t}^i, a_{t:t+\tau_F}^i, g_t^i)$ from dataset \mathcal{D} , where i is the sample number, $I_{t-\tau_P:t}^i$ is a sequence of past visual observations starting from $t - \tau_P$ and ending at t , $a_{t:t+\tau_F}^i$ is a sequence of future actions starting from t and ending at $t + \tau_F$, and g_t^i is the current goal at time t instantiated as an image in the future $I_{t+\tau_F}^i$. τ_F is the number of frames in the future and τ_P is the number of frames in the past. We then feed $I_{t-\tau_P:t}^i$ to p_ϕ , typically a CNN, and all the embeddings to a transformer encoder [46], as well as $a_{t:t+\tau_F}^i$ to f_ξ as part of another transformer encoder, to learn image Z^i and action Z^a embeddings, respectively. Each transformer contains an additional context token to capture the continuous information among frames. We feed g_t^i to p_ϕ to generate goal embedding Z^g . Finally, we use VANP’s objective function to learn ϕ and ξ :

$$\begin{aligned} \mathcal{L}_{\text{VANP}}(Z^i, Z^g, Z^a) &= \lambda \mathcal{L}_{\text{VICReg}}(Z^i, Z^g) \\ &\quad + (1 - \lambda) \mathcal{L}_{\text{VICReg}}(Z^i, Z^a), \end{aligned} \quad (1)$$

where λ is the importance of each term, and $\mathcal{L}_{\text{VICReg}}$ is the VICReg objective function [23] defined as:

$$\begin{aligned} \mathcal{L}_{\text{VICReg}}(Z^1, Z^2) &= \mu^1 s(Z^1, Z^2) \\ &\quad + \mu^2 [v(Z^1) + v(Z^2)] \\ &\quad + \mu^3 [c(Z^1) + c(Z^2)]. \end{aligned} \quad (2)$$

s is the distance between embedding spaces, v and c are the variance and covariance of each embedding respectively. μ^1 ,

μ^2 , and μ^3 are hyper-parameters controlling the effectiveness of each term. Leveraging VICReg’s objective function offers the advantage of circumventing the need for negative samples, which, as mentioned above, is challenging to define within the action space for navigation tasks. We also compare VICReg’s performance against BarlowTwins used by Nazeri *et al.* [68] and observe that BarlowTwins tends to prioritize redundant scene features over those with greater relevance to navigation resulting in degraded performance.

C. Implementation Details

We implement VANP with PyTorch [69] and the training is performed on a single A5000 GPU with 24GB memory¹.

Model architecture: Considering the limited computation resources onboard most mobile robots, we choose ResNet-50 [24] without the classification head as a low-latency image encoder for p_ϕ and we call it VANP-50. We use two TransformerEncoders with additional context vectors [9], [21], [70] with four layers and four heads as the final image and action encoders to produce the embeddings of $Z^i, Z^a \in \mathbb{R}^{512}$. Both encoders are followed by MLPs with three layers as the projection heads to generate the final $Z^i, Z^a \in \mathbb{R}^{1024}$. We apply the same p_ϕ to the goal image to generate $Z^g \in \mathbb{R}^{512}$. A critical challenge arises from the inherent differences in modalities between the two networks generating the embeddings, leading to significant variations in their output ranges. To address this discrepancy and ensure effective integration, we initialize all deep networks using the Kaiming Normal initialization [24] with a mean of zero and a variance of one. In the context of the downstream model, an MLP is appended to the Transformer’s context vector to predict trajectories at three and five seconds into the future, enabling the evaluation of how the extracted features influence both short-term and long-term interactions.

Optimization: We use the ADAMW optimizer [71] and train the model for 200 epochs with a batch size of 2048 and a learning rate of $5e^{-4}$. We observe that large batch sizes add more variation to the update stage and improve learning. To ensure a fair comparison, all models are trained for 50 epochs using the same optimizer and hyperparameters during downstream training. The sole exception is the end-to-end model, which requires 100 epochs to guarantee convergence.

Dataset: We leverage a selection of two unique datasets: SCAND [25] and MuSoHu [72], both of which encapsulate robot and human navigation data from the egocentric perspective. Both real-world datasets are collected in a variety of natural crowded public spaces. MuSoHu comprises approximately 20 hours of data captured from human egocentric motion. The recordings capture human walking patterns in public spaces, providing insights for learning human-like, socially compliant navigation behaviors. SCAND is an autonomous robot navigation dataset that captures 8.7 hours of human-teleoperated robot navigation demonstrations in naturally crowded public spaces on a university campus. A fundamental limitation of SSL models is their susceptibility

¹<https://github.com/mhnazeri/VANP>

TABLE I: **Downstream Performance.** Comparison of the performance of the visual encoders with different pre-training methods on unseen data. Models denoted by an ⌚ require double the training time compared to models with ⌚

Type	Method	Weight	Single-frame	Multiple-frame	Frozen ⌚		Fine-tuned ⌚	
					3s	5s	3s	5s
End-to-End ⌚	Resnet-50	Random	✓	✗	-	-	0.116	0.307
	ResnetTransformer	Random	✗	✓	-	-	0.113	0.320
Backbone Supervised ⌚	Resnet-50	ImageNet	✓	✗	0.129	0.356	0.129	0.342
	ResnetTransformer	ImageNet	✗	✓	0.169	0.435	0.107	0.292
Backbone Self-Supervised ⌚	Resnet-50	VANP	✓	✗	0.144	0.374	0.103	0.272
	ResnetTransformer	VANP	✗	✓	0.133	0.342	0.114	0.319

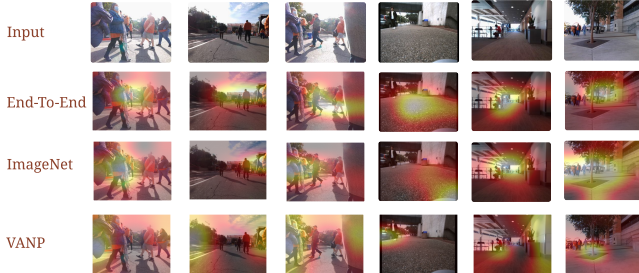


Fig. 3: **Qualitative Comparison.** Comparison of the last layer activation maps among different methods on unseen scenarios.

to data quality [20]. As we will discuss in the limitations section (Sec. IV-D), VANP is similarly affected, particularly in scenarios where there is no change in a sequence of images as shown in Fig. 4. To minimize data ambiguity and noise, a subset of the two datasets are carefully curated, ensuring representation of both indoor and outdoor scenes. The resulting dataset, comprising approximately 11,000 samples, was used for both pre-training and training phases. Additionally, a separate set of 8,000 unseen samples are used for downstream navigation task evaluation. For pretext task training, we set τ_P and τ_F to 6 and 20 respectively and use a sequence of images $I_{t-\tau_P+1:t} \in \mathbb{R}^{\tau_P \times 98 \times 126}$ along with a goal image $g_t \in \mathbb{R}^{98 \times 126}$ and a sequence of actions $a_{t:t+\tau_F-1} \in \mathbb{R}^{\tau_F \times 2}$ parsed at 4 Hz, comprising of 1.5 seconds in the past and 5 seconds in the future. For the downstream task, we use a sequence of past observations $I_{t-\tau_P+1:t} \in \mathbb{R}^{\tau_P \times 98 \times 126}$ along with the polar coordinates of the next local goal $g \in \mathbb{R}^2$ parsed at 4 Hz, containing 1.5 seconds history as the network input to produce the actions $\mathcal{A}_{t:t+\tau_F-1} \in \mathbb{R}^{\tau_F \times 2}$ for three and five seconds in the future.

IV. EXPERIMENTAL RESULTS

We provide experimental results using VANP compared against a ResNet-50 pre-trained on ImageNet and end-to-end from scratch as baselines.

A. Results Discussion

We assess the efficacy of VANP pretext training by quantitatively comparing its performance with that of a ResNet-50 model [24] pre-trained on the ImageNet ILSVRC-2012 dataset [12]. This serves as the baseline alongside

another ResNet-50 model trained end-to-end with randomly initialized weights. To guarantee a fair comparison, the architectures of all other components within the downstream task remain unchanged. Table I presents the mean squared error between the predicted and ground truth trajectories for short-(three seconds) and long-term (five seconds) interactions under two conditions. In the first condition, only the goal encoder and controller are trained during the downstream navigation task, while the image encoder weights are frozen. In the second condition, we compare the performance by unfreezing the image encoder weights to enable fine-tuning.

The results in Table I demonstrate that VANP achieves comparable performance to the end-to-end trained model while requiring only half the training time. Furthermore, VANP pre-trained model achieves comparable performance to ImageNet model with only 0.08% of the data size required by ImageNet, highlighting how informative the extracted representations are for navigation.

When provided with a sequence of past observations, VANP exhibits a superior ability (0.342) to utilize this additional data compared to ImageNet model when frozen (0.435). Although the ImageNet weights appear unable to leverage the temporal features provided by the transformer component when freezing its weights (Table I, row four compared against row three), fine-tuning the ImageNet model leads to performance improvement from 0.435 to 0.292, suggesting that it can better capture underlying temporal features provided by the Transformer through fine-tuning.

However, we do not see such an improvement in the case of VANP. The negligible improvement in accuracy from 0.342 to 0.319 for VANP during fine-tuning can be attributed to two reasons. First, the focus on multiple navigation-related visual regions of VANP’s pre-trained weights (Fig. 3 last row) impedes adaptation/forgetting during fine-tuning compared to the ImageNet weights. Second, the temporal features from the Transformer are already in VANP weights and therefore does not require much fine-tuning. Overall, it is likely that forgetting/updating weights can be easier when the visual encoder is trained using only one single scalar instructive feedback (i.e., training loss) rather than pre-trained on richer instructive signals, i.e., VANP’s pre-training objective signal.

Interestingly, during frozen evaluation with only one image as input, the frozen pre-trained ImageNet model (Table I,

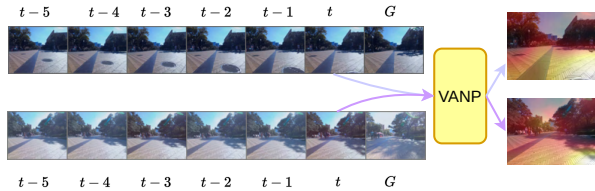


Fig. 4: **Failure Cases.** Samples without any important intra-frame changes cause the model to collapse.

TABLE II: **Ablations.** Ablation study on the role of each module on the downstream navigation task performance.

Information	3s	5s
Actions	0.167	0.499
Goal	0.160	0.392
Actions+GoalIn	0.155	0.386
Actions+GoalOut	0.144	0.383
Augmentations	0.133	0.342

row three) achieves the best performance. This finding warrants further investigation. One assumption is that in test cases, the salient object has a stronger influence on the trajectory and aligns better with the single, scalar form of instructive feedback provided. However, during fine-tuning, it is clear that the single image does not outperform models utilizing Transformer temporal features (0.342) while the VANP model benefits from these features even with only one image as input (0.272).

Visual inspection of the learned activation maps on the last layer of ResNet-50 (Fig. 3) reveals distinct characteristics across the models. The last row on Fig. 3 shows that the VANP pre-trained model exhibits activation maps with a higher degree of relevance to navigation tasks, focusing on features such as paths and obstacles while the ImageNet pre-trained model (Fig. 3 third row) primarily focuses on salient objects within the environment, which might not be directly related to navigation. Another difference between VANP and the end-to-end model (Fig. 3 second row) is that the end-to-end model tends to concentrate on a single critical region significantly impacting the trajectory, likely due to its limited instructive signal during training, i.e., minimizing the distance between predicted and ground truth trajectories. Conversely, VANP demonstrates the ability to extract information from multiple regions, potentially benefiting from the richer information provided by the goal image and future actions during the pre-training stage. However, as mentioned above, this richness impedes adaptation during fine-tuning.

We observe instances where the attention of all models shifts to seemingly irrelevant aspects. In the case of VANP, we posit that this may be due to the robot’s sharp turns temporarily obscuring the goal image from the current frame.

B. Ablations

To investigate the most effective approach for correlating visual and action spaces, we conduct a series of ablation studies, in which we report the mean squared distance of the

predicted trajectory from the ground truth in three and five seconds in the future in Table II.

Role of Different Training Signals: We assessed the individual contributions of various self-supervised training signals by changing the value of λ between 0 and 1 in Eq. 1. Our findings reveal that while action signals provide valuable navigational cues, their sparsity often hinders their effectiveness in downstream navigation tasks, especially during long-term interactions. Conversely, information derived from the goal, while occasionally exhibiting redundancy, improved performance from 0.499 to 0.392 during long-term interactions over using only actions due to informative cues alongside the redundant elements. However, this redundancy poses challenges for the policy network, which can be remedied by more training epochs and a deeper policy network. By combining these two embeddings as the self-supervision signal, the final model can effectively learn informative features while mitigating the impact of redundant information within the embedding.

Leveraging Goal Information: We further investigated the optimal utilization of future goal information. Our findings suggest that employing the goal solely as a supervision signal (shown as Actions+GoalOut in Table II) proves more effective in facilitating the model’s learning of visual features compared to incorporating the goal directly within the Transformer architecture (shown as Actions+GoalIn in Table II). The Transformer’s ability to capture temporal changes from the current to the goal frame is only helpful when the goal is visible from the current frame.

Augmentations: Data augmentation is a standard technique employed to enhance model generalization by introducing variability into the dataset. We follow the augmentation scheme outlined by Bardes *et al.* [23] and the result is shown as Augmentations in Table II. We observe that random cropping is particularly critical for VANP, especially in scenarios exemplified by Fig. 4, as it introduces inter-frame variation. This augmentation strategy relaxes the assumption of carefully curated data and enables an expansion of the dataset from 11,000 to 26,042 samples to include even ambiguous and noisy samples with a little performance hit.

C. Robot Deployment

To demonstrate the practical applicability of the learned visual features for navigation, a proof-of-concept demonstration of VANP-18 with a moving goal objective [8] is deployed on a Clearpath Jackal robot. The obstacle avoidance capabilities of VANP are evaluated under controlled conditions. In these experiments, a static obstacle is initially positioned in the robot’s path. Subsequent trials involve a dynamic obstacle, simulated by a human pedestrian. Results indicate that VANP exhibits an ability to detect and avoid both static and dynamic obstructions in the majority of test cases. It is important to note that while VANP demonstrates capabilities in object avoidance, it encounters difficulties in navigating around minor obstacles, a limitation likely attributable to restricted visibility conditions. The supplement-

tary video provides a record of these experiments². Despite VANP’s intended versatility across diverse environmental conditions, inherent limitations considering safety only allow it to work in uncluttered environments, as elaborated in the subsequent section.

D. Limitations

We identify multiple key limitations of the VANP pre-training approach. First, our analysis of the learned kernels suggests that VANP performs more effectively when the goal image is directly visible from the current image, likely due to its reliance on image correlation for learning. While this is helpful for Visual-Goal navigation task, it highlights a potential limitation in generalizability to scenarios where the goal location may not be directly visible from the starting point. Second, in large-scale datasets likely with a significant amount of noise, scaling VANP poses a potential challenge, considering its need for high-quality self-supervision during pre-training can result in many changes in learned activation maps between epochs. As can be seen in Fig. 4, the VANP objective is unable to learn from scenarios where there is no intra-frame change as the time passes. This limitation can be alleviated with augmentations, particularly random cropping, but it does not eliminate it. Additionally, our current findings are based on a static dataset and may not directly translate to challenging real-world navigation tasks that involve dynamic environments and unforeseen obstacles. Further research is needed to evaluate VANP’s performance in these more complex scenarios.

V. CONCLUSIONS AND FUTURE WORK

In this work, we propose a self-supervised learning approach to train visual encoder models specifically designed for visual navigation. This approach is motivated by the observation that humans only pay attention to specific navigation-relevant regions of their frontal view to efficiently make navigation decisions. By reversing this observation, we use the navigation decisions to extract only visual features that are relevant to the navigation task, unlike computer vision models that mainly extract salient details, which are potentially irrelevant to navigation tasks and can therefore lead to confusion for neural-based controllers. To achieve this, we leverage two Transformer Encoders to embed past visual observation, future actions, and a goal image, then we maximize the information between these embeddings using VANP’s objective function to learn visual backbone weights.

Furthermore, the VANP objective function facilitates the integration of additional embeddings derived from diverse modalities, including depth data and semantic information or inputs from other sensors such as LiDARs [73]. Studying the effectiveness of this enrichment of the embedding space with supplementary information for downstream navigation tasks can be a potential future work. Another future direction is to merge datasets from different environments, such as indoor [74], [75], outdoor [76], [77], off-road [78]–[80], and

social environments [81], [82], to extend the generalizability of the proposed VANP approach. More real-world experiments can support all these future directions and scale up the model to larger datasets.

REFERENCES

- [1] D. A. Pomerleau, “ALVINN: An autonomous land vehicle in a neural network,” in *Advances in Neural Information Processing Systems*, D. Touretzky, Ed., vol. 1. Morgan-Kaufmann, 1988.
- [2] M. Bojarski, D. Del Testa, D. Dworakowski, B. Firner, B. Flepp, P. Goyal, L. D. Jackel, M. Monfort, U. Muller, J. Zhang, X. Zhang, J. Zhao, and K. Zieba, “End to End Learning for Self-Driving Cars,” Apr. 2016.
- [3] F. Codevilla, M. Muller, A. Lopez, V. Koltun, and A. Dosovitskiy, “End-to-End Driving Via Conditional Imitation Learning,” in *2018 IEEE International Conference on Robotics and Automation (ICRA)*. Brisbane, QLD: IEEE, May 2018, pp. 4693–4700.
- [4] F. Codevilla, E. Santana, A. M. Lopez, and A. Gaidon, “Exploring the Limitations of Behavior Cloning for Autonomous Driving,” in *The IEEE International Conference on Computer Vision (ICCV)*. IEEE, Oct. 2019.
- [5] E. Ohn-Bar, A. Prakash, A. Behl, K. Chitta, and A. Geiger, “Learning Situational Driving,” in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Jun. 2020.
- [6] D. Chen, B. Zhou, V. Koltun, and P. Krähenbühl, “Learning by cheating,” in *Proceedings of the Conference on Robot Learning*, ser. Proceedings of Machine Learning Research, L. P. Kaelbling, D. Kragic, and K. Sugiura, Eds., vol. 100. PMLR, 2020-11-30, 2020, pp. 66–75.
- [7] M. H. Nazeri and M. Bohlouli, “Exploring Reflective Limitation of Behavior Cloning in Autonomous Vehicles,” in *2021 IEEE International Conference on Data Mining (ICDM)*. Auckland, New Zealand: IEEE, Dec. 2021, pp. 1252–1257.
- [8] X. Xiao, B. Liu, G. Warnell, and P. Stone, “Motion planning and control for mobile robot navigation using machine learning: a survey,” *Autonomous Robots*, vol. 46, no. 5, pp. 569–597, 2022.
- [9] 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,” Jun. 2021.
- [10] D. Shah, A. Sridhar, A. Bhorkar, N. Hirose, and S. Levine, “GNM: A general navigation model to drive any robot,” in *International Conference on Robotics and Automation (ICRA)*, 2023.
- [11] D. Shah, A. Sridhar, N. Dashora, K. Stachowicz, K. Black, N. Hirose, and S. Levine, “ViNT: A foundation model for visual navigation,” in *7th Annual Conference on Robot Learning*, 2023.
- [12] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “ImageNet: A large-scale hierarchical image database,” in *2009 IEEE Conference on Computer Vision and Pattern Recognition*, Jun. 2009, pp. 248–255.
- [13] K. Vishniakov, Z. Shen, and Z. Liu, “ConvNet vs Transformer, Supervised vs CLIP: Beyond ImageNet Accuracy,” Jan. 2024.
- [14] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, “A simple framework for contrastive learning of visual representations,” in *International Conference on Machine Learning*. PMLR, 2020, pp. 1597–1607.
- [15] J. Zbontar, L. Jing, I. Misra, Y. LeCun, and S. Deny, “Barlow Twins: Self-Supervised Learning via Redundancy Reduction,” Jun. 2021.
- [16] M. Oquab, T. Darcet, T. Moutakanni, H. Vo, M. Szafraniec, V. Khalidov, P. Fernandez, D. Haziza, F. Massa, A. El-Nouby, M. Assran, N. Ballas, W. Galuba, R. Howes, P.-Y. Huang, S.-W. Li, I. Misra, M. Rabbat, V. Sharma, G. Synnaeve, H. Xu, H. Jegou, J. Mairal, P. Labatut, A. Joulin, and P. Bojanowski, “DINOv2: Learning Robust Visual Features without Supervision,” Apr. 2023.
- [17] A. Bardes, J. Ponce, and Y. LeCun, “MC-JEPA: A Joint-Embedding Predictive Architecture for Self-Supervised Learning of Motion and Content Features,” Jul. 2023.
- [18] S. Gidaris, P. Singh, and N. Komodakis, “Unsupervised representation learning by predicting image rotations,” *arXiv preprint arXiv:1803.07728*, 2018.
- [19] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick, “Masked autoencoders are scalable vision learners,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 16 000–16 009.

²<https://youtu.be/SEuD9hkwXxQ>

- [20] R. Balestriero, M. Ibrahim, V. Sobal, A. Morcos, S. Shekhar, T. Goldstein, F. Bordes, A. Bardes, G. Mialon, Y. Tian, A. Schwarzschild, A. G. Wilson, J. Geiping, Q. Garrido, P. Fernandez, A. Bar, H. Pirsiavash, Y. LeCun, and M. Goldblum, "A Cookbook of Self-Supervised Learning," Apr. 2023.
- [21] 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)*, J. Burstein, C. Doran, and T. Solorio, Eds. Minneapolis, Minnesota: Association for Computational Linguistics, Jun. 2019, pp. 4171–4186.
- [22] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun, "CARLA: An open urban driving simulator," in *Proceedings of the 1st Annual Conference on Robot Learning*, 2017, pp. 1–16.
- [23] A. Bardes, J. Ponce, and Y. LeCun, "VICReg: Variance-invariance-covariance regularization for self-supervised learning," in *International Conference on Learning Representations*, 2022.
- [24] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, 2015.
- [25] H. Karnan, A. Nair, X. Xiao, G. Warnell, S. Pirk, A. Toshev, J. Hart, J. Biswas, and P. Stone, "Socially compliant navigation dataset (scand): A large-scale dataset of demonstrations for social navigation," *IEEE Robotics and Automation Letters*, vol. 7, no. 4, pp. 11 807–11 814, 2022.
- [26] J. Bromley, I. Guyon, Y. LeCun, E. Säckinger, and R. Shah, "Signature verification using a "Siamese" time delay neural network," in *Advances in Neural Information Processing Systems*, J. Cowan, G. Tesauoro, and J. Alspector, Eds., vol. 6. Morgan-Kaufmann, 1993.
- [27] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9729–9738.
- [28] K. S. Sikand, S. Rabiee, A. Uccello, X. Xiao, G. Warnell, and J. Biswas, "Visual representation learning for preference-aware path planning," in *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022, pp. 11 303–11 309.
- [29] J.-B. Grill, F. Strub, F. Altché, C. Tallec, P. Richemond, E. Buchatskaya, C. Doersch, B. Avila Pires, Z. Guo, M. Gheshlaghi Azar *et al.*, "Bootstrap your own latent—a new approach to self-supervised learning," *Advances in neural information processing systems*, vol. 33, pp. 21 271–21 284, 2020.
- [30] X. Chen and K. He, "Exploring Simple Siamese Representation Learning," Nov. 2020.
- [31] D. Chen, V. Koltun, and P. Krähenbühl, "Learning to drive from a world on rails," in *ICCV*, 2021.
- [32] B. Jaeger and A. Geiger, "An Invitation to Deep Reinforcement Learning," Dec. 2023.
- [33] Y. Bengio, A. Courville, and P. Vincent, "Representation Learning: A Review and New Perspectives," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 8, pp. 1798–1828, Aug. 2013, conference Name: IEEE Transactions on Pattern Analysis and Machine Intelligence.
- [34] A. Payandeh, K. T. Baghaei, P. Fayyazsanavi, S. B. Ramezani, Z. Chen, and S. Rahimi, "Deep representation learning: Fundamentals, technologies, applications, and open challenges," *IEEE Access*, vol. 11, pp. 137 621–137 659, 2023.
- [35] R. Dadashi, L. Hussenot, D. Vincent, S. Girgin, A. Raichuk, M. Geist, and O. Pietquin, "Continuous control with action quantization from demonstrations," in *International Conference on Machine Learning*, PMLR, 2022, pp. 4537–4557.
- [36] J. Luo, P. Dong, Y. Zhai, Y. Ma, and S. Levine, "RLIF: Interactive Imitation Learning as Reinforcement Learning," Nov. 2023.
- [37] W. Huang, Y. Zhou, X. He, and C. Lv, "Goal-guided Transformer-enabled Reinforcement Learning for Efficient Autonomous Navigation," Jan. 2023.
- [38] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan, K. Hausman, A. Herzog, J. Hsu, J. Ibarz, B. Ichter, A. Irpan, T. Jackson, S. Jesmonth, N. J. Joshi, R. Julian, D. Kalashnikov, Y. Kuang, I. Leal, K.-H. Lee, S. Levine, Y. Lu, U. Malla, D. Manjunath, I. Mordatch, O. Nachum, C. Parada, J. Peralta, E. Perez, K. Pertsch, J. Quiambao, K. Rao, M. Ryoo, G. Salazar, P. Sanketi, K. Sayed, J. Singh, S. Sontakke, A. Stone, C. Tan, H. Tran, V. Vanhoucke, S. Vega, Q. Vuong, F. Xia, T. Xiao, P. Xu, S. Xu, T. Yu, and B. Zitkovich, "RT-1: Robotics Transformer for Real-World Control at Scale," Aug. 2023.
- [39] N. Di Palo, A. Byravan, L. Hasenclever, M. Wulfmeier, N. Heess, and M. Riedmiller, "Towards A Unified Agent with Foundation Models," Jul. 2023.
- [40] A. Hiranaka, M. Hwang, S. Lee, C. Wang, L. Fei-Fei, J. Wu, and R. Zhang, "Primitive Skill-based Robot Learning from Human Evaluative Feedback," Aug. 2023.
- [41] S. Nair, A. Rajeswaran, V. Kumar, C. Finn, and A. Gupta, "R3M: A universal visual representation for robot manipulation," in *Conference on Robot Learning*. PMLR, 2023, pp. 892–909.
- [42] K. Grauman, A. Westbury, E. Byrne, Z. Chavis, A. Furnari, R. Girdhar, J. Hamburger, H. Jiang, M. Liu, X. Liu, M. Martin, T. Nagarajan, I. Radosavovic, S. K. Ramakrishnan, F. Ryan, J. Sharma, M. Wray, M. Xu, E. Z. Xu, C. Zhao, S. Bansal, D. Batra, V. Cartillier, S. Crane, T. Do, M. Doulaty, A. Erapalli, C. Feichtenhofer, A. Fragomeni, Q. Fu, C. Fuegen, A. Gebreselasie, C. Gonzalez, J. Hillis, X. Huang, Y. Huang, W. Jia, W. Khoo, J. Kolar, S. Kottur, A. Kumar, F. Landini, C. Li, Y. Li, Z. Li, K. Mangalam, R. Modhugu, J. Munro, T. Murrell, T. Nishiyasu, W. Price, P. R. Puentes, M. Ramazanova, L. Sari, K. Somasundaram, A. Southerland, Y. Sugano, R. Tao, M. Vo, Y. Wang, X. Wu, T. Yagi, Y. Zhu, P. Arbelaez, D. Crandall, D. Damen, G. M. Farinella, B. Ghanem, V. K. Ithapu, C. V. Jawahar, H. Joo, K. Kitani, H. Li, R. Newcombe, A. Oliva, H. S. Park, J. M. Rehg, Y. Sato, J. Shi, M. Z. Shou, A. Torralba, L. Torresani, M. Yan, and J. Malik, "Ego4D: Around the World in 3,000 Hours of Egocentric Video," in *IEEE/CVF Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [43] M. Shridhar, L. Manuelli, and D. Fox, "Cliport: What and where pathways for robotic manipulation," in *Conference on Robot Learning*. PMLR, 2022, pp. 894–906.
- [44] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever, "Learning transferable visual models from natural language supervision," in *Proceedings of the 38th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, M. Meila and T. Zhang, Eds., vol. 139. PMLR, Jul. 2021, pp. 8748–8763.
- [45] A. van den Oord, O. Vinyals, and K. Kavukcuoglu, "Neural discrete representation learning," *arXiv preprint arXiv:1711.00937*, 2017.
- [46] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [47] A. R. Zamir, A. Sax, W. Shen, L. J. Guibas, J. Malik, and S. Savarese, "Taskonomy: Disentangling task transfer learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2018.
- [48] W. Shen, D. Xu, Y. Zhu, L. Fei-Fei, L. Guibas, and S. Savarese, "Situational Fusion of Visual Representation for Visual Navigation," in *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*. Seoul, Korea (South): IEEE, Oct. 2019, pp. 2881–2890.
- [49] H.-K. Yang, T.-C. Chiang, T.-R. Liu, C.-W. Huang, J.-M. Liu, and C.-Y. Lee, "Virtual Guidance as a Mid-level Representation for Navigation," Sep. 2023.
- [50] H. Karnan, E. Yang, D. Farkash, G. Warnell, J. Biswas, and P. Stone, "STERLING: Self-Supervised Terrain Representation Learning from Unconstrained Robot Experience," Oct. 2023.
- [51] A. Pokhrel, A. Datar, M. Nazeri, and X. Xiao, "CAHSOR: Competence-Aware High-Speed Off-Road Ground Navigation in SE(3)," *arXiv preprint arXiv:2402.07065*, 2024.
- [52] A. Eftekhari, K.-H. Zeng, J. Duan, A. Farhadi, A. Kembhavi, and R. Krishna, "Selective Visual Representations Improve Convergence and Generalization for Embodied AI," Nov. 2023.
- [53] Y. Wang, C.-Y. Ko, and P. Agrawal, "Visual pre-training for navigation: What can we learn from noise?" in *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2023, pp. 3897–3902.
- [54] H. Wang, A. H. Tan, and G. Nejat, "NavFormer: A Transformer Architecture for Robot Target-Driven Navigation in Unknown and Dynamic Environments," *arXiv preprint arXiv:2402.06838*, 2024.
- [55] M. Ahn, A. Brohan, N. Brown, Y. Chebotar, O. Cortes, B. David, C. Finn, C. Fu, K. Gopalakrishnan, K. Hausman *et al.*, "Do as i can, not as i say: Grounding language in robotic affordances," *arXiv preprint arXiv:2204.01691*, 2022.
- [56] S. Li, X. Puig, C. Paxton, Y. Du, C. Wang, L. Fan, T. Chen, D.-A. Huang, E. Akyürek, A. Anandkumar *et al.*, "Pre-trained language mod-

- els for interactive decision-making,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 31 199–31 212, 2022.
- [57] M. Reid, Y. Yamada, and S. S. Gu, “Can wikipedia help offline reinforcement learning?” *arXiv preprint arXiv:2201.12122*, 2022.
- [58] B. Baker, I. Akkaya, P. Zhokov, J. Huizinga, J. Tang, A. Ecoffet, B. Houghton, R. Sampedro, and J. Clune, “Video pretraining (vpt): Learning to act by watching unlabeled online videos,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 24 639–24 654, 2022.
- [59] D. Shah, A. Sridhar, A. Bhorkar, N. Hirose, and S. Levine, “GNM: A General Navigation Model to Drive Any Robot,” Oct. 2022.
- [60] V. Shah, F. Träuble, A. Malik, H. Larochelle, M. Mozer, S. Arora, Y. Bengio, and A. Goyal, “Unlearning via Sparse Representations,” Nov. 2023.
- [61] D. Fox, W. Burgard, and S. Thrun, “The dynamic window approach to collision avoidance,” *IEEE Robotics & Automation Magazine*, vol. 4, no. 1, pp. 23–33, 1997.
- [62] S. Quinlan and O. Khatib, “Elastic bands: Connecting path planning and control,” in *[1993] Proceedings IEEE International Conference on Robotics and Automation*. IEEE, 1993, pp. 802–807.
- [63] X. Xiao, B. Liu, G. Warnell, and P. Stone, “Toward agile maneuvers in highly constrained spaces: Learning from hallucination,” *IEEE Robotics and Automation Letters*, vol. 6, no. 2, pp. 1503–1510, 2021.
- [64] S. Ravi, G. Wang, S. Satewar, X. Xiao, G. Warnell, J. Biswas, and P. Stone, “Visually adaptive geometric navigation,” in *2023 IEEE International Symposium on Safety, Security, and Rescue Robotics (SSRR)*. IEEE, 2023.
- [65] X. Xiao, Z. Xu, Z. Wang, Y. Song, G. Warnell, P. Stone, T. Zhang, S. Ravi, G. Wang, H. Karnan *et al.*, “Autonomous ground navigation in highly constrained spaces: Lessons learned from the benchmark autonomous robot navigation challenge at icra 2022 [competitions],” *IEEE Robotics & Automation Magazine*, vol. 29, no. 4, pp. 148–156, 2022.
- [66] X. Xiao, Z. Wang, Z. Xu, B. Liu, G. Warnell, G. Dhamankar, A. Nair, and P. Stone, “Appl: Adaptive planner parameter learning,” *Robotics and Autonomous Systems*, vol. 154, p. 104132, 2022.
- [67] T. Chen, S. Kornblith, K. Swersky, M. Norouzi, and G. Hinton, “Big self-supervised models are strong semi-supervised learners,” *arXiv preprint arXiv:2006.10029*, 2020.
- [68] M. Nazeri, J. Wang, A. payandeh, and X. Xiao, “VANP: Self-supervised vision-action pretraining for navigation,” in *Bridging the Gap between Cognitive Science and Robot Learning in the Real World: Progresses and New Directions*, 2024. [Online]. Available: <https://openreview.net/forum?id=MOI3CESxR2>
- [69] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimeshein, 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 *Advances in Neural Information Processing Systems 32*. Curran Associates, Inc., 2019, pp. 8024–8035.
- [70] T. Darcet, M. Oquab, J. Mairal, and P. Bojanowski, “Vision Transformers Need Registers,” Sep. 2023.
- [71] I. Loshchilov and F. Hutter, “Decoupled Weight Decay Regularization,” Jan. 2019.
- [72] D. M. Nguyen, M. Nazeri, A. Payandeh, A. Datar, and X. Xiao, “Toward human-like social robot navigation: A large-scale, multi-modal, social human navigation dataset,” in *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2023.
- [73] B. Panigrahi, A. H. Raj, M. Nazeri, and X. Xiao, “A study on learning social robot navigation with multimodal perception,” *arXiv preprint arXiv:2309.12568*, 2023.
- [74] B. Liu, X. Xiao, and P. Stone, “A lifelong learning approach to mobile robot navigation,” *IEEE Robotics and Automation Letters*, vol. 6, no. 2, pp. 1090–1096, 2021.
- [75] P. Atreya, H. Karnan, K. S. Sikand, X. Xiao, S. Rabiee, and J. Biswas, “High-speed accurate robot control using learned forward kinodynamics and non-linear least squares optimization,” in *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2022, pp. 11 789–11 795.
- [76] X. Xiao, J. Biswas, and P. Stone, “Learning inverse kinodynamics for accurate high-speed off-road navigation on unstructured terrain,” *IEEE Robotics and Automation Letters*, vol. 6, no. 3, pp. 6054–6060, 2021.
- [77] H. Karnan, K. S. Sikand, P. Atreya, S. Rabiee, X. Xiao, G. Warnell, P. Stone, and J. Biswas, “Vi-ikd: High-speed accurate off-road navigation using learned visual-inertial inverse kinodynamics,” in *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2022, pp. 3294–3301.
- [78] A. Datar, C. Pan, M. Nazeri, and X. Xiao, “Toward wheeled mobility on vertically challenging terrain: Platforms, datasets, and algorithms,” in *2024 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2024.
- [79] A. Datar, C. Pan, and X. Xiao, “Learning to model and plan for wheeled mobility on vertically challenging terrain,” *arXiv preprint arXiv:2306.11611*, 2023.
- [80] A. Datar, C. Pan, M. Nazeri, A. Pokhrel, and X. Xiao, “Terrain-attentive learning for efficient 6-dof kinodynamic modeling on vertically challenging terrain,” *arXiv preprint arXiv:2403.16419*, 2024.
- [81] X. Xiao, T. Zhang, K. M. Choromanski, T.-W. E. Lee, A. Francis, J. Varley, S. Tu, S. Singh, P. Xu, F. Xia, S. M. Persson, L. Takayama, R. Frostig, J. Tan, C. Parada, and V. Sindhwani, “Learning model predictive controllers with real-time attention for real-world navigation,” in *Conference on robot learning*. PMLR, 2022.
- [82] R. Mirsky, X. Xiao, J. Hart, and P. Stone, “Conflict avoidance in social navigation—a survey,” *ACM Transactions on Human-Robot Interaction*, vol. 13, no. 1, pp. 1–36, 2024.