

# VANP: Self-Supervised Vision-Action Pretraining for Navigation

Anonymous Author(s)

Affiliation

Address

email

1       **Abstract:** Self-supervised learning has revolutionized the fields of computer vi-  
2 sion and natural language processing. Despite its potential, its application to  
3 robotic navigation tasks remains under-explored. This is due to the difficulty of  
4 defining effective self-supervision signals for robotics. Fortunately, with the re-  
5 cent development of many large-scale robotic navigation datasets with a variety of  
6 sensor and action data that can be used as self-supervision signals, self-supervised  
7 learning has become a viable approach for robotic navigation tasks. In this work,  
8 we propose a self-supervised method for learning visual features for end-to-end  
9 robotic navigation systems, using actions as the supervisory signal. This approach  
10 is motivated by the observation that humans tend to focus on specific regions of  
11 their frontal view in order to make navigation decisions and produce navigation  
12 actions. We reverse this procedure, using future actions to learn only the visual  
13 features that are important for navigation, as opposed to extracted features by  
14 conventional computer vision models that tend to extract every detail of the envi-  
15 ronment that can be misleading to a downstream navigation controller. Our results  
16 show that this approach enables small convolutional neural network-based visual  
17 encoders to achieve performance comparable to large vision foundation models  
18 trained on billions of images. This demonstrates the scalability and effectiveness  
19 of our self-supervised learning method for robotic navigation.

20       **Keywords:** Self-supervised, Navigation, Learning

## 21 1 Introduction

22 Recent advances in computer vision and deep learning rely on the development of increasingly large  
23 and complex Deep Neural Networks (DNNs) [1, 2, 3, 4]. However, training these DNNs from scratch  
24 can be computationally expensive and requires a large amount of computing resources [5, 6, 7, 8, 9].  
25 Self-Supervised Learning (SSL) [10, 11, 12, 13] is a machine learning paradigm that can mitigate  
26 the need for annotated data by enabling DNNs to first pre-train the model from unlabeled data  
27 and then quickly fine-tune to adapt to specific tasks, avoiding the need to re-train everything from  
28 scratch, i.e., SSL trains DNNs to complete a pretext task that does not require labels. For example,  
29 DNNs might be trained to predict the rotation of an image [14] or to reconstruct an image from  
30 its corrupted/obstructed version [15]. By completing these pretext tasks, DNNs learn to extract  
31 meaningful features from the data, which can then be used to solve downstream tasks such as image  
32 classification and object detection [16].

33 Despite the effectiveness of SSL in a variety of computer vision tasks, challenges still remain that  
34 need to be addressed before SSL can be widely adopted in robotics applications. For example,  
35 SSL models typically require a large amount of data to train, which can be difficult to obtain in  
36 robotics settings. Additionally, SSL models trained on computer vision datasets such as ImageNet  
37 ILSVRC-2012 [17] and COCO [18] may not generalize well to robotic navigation tasks, which  
38 contain a significant amount of dynamic scenes of moving agents that can affect a robot’s trajectory.



Figure 1: An illustrative example of the difference between attentive regions from computer vision models (left) and navigation models (right). With the assumption that everyone respects social etiquette, some regions become redundant when making navigational decisions.

39 As shown in Figure 1 left, complex DNN models trained on vision tasks may extract all the existing  
 40 features in the scene that may confuse downstream navigation controllers based on Neural Networks  
 41 (NNs) by providing too much information and leading to improper actions. Although these features  
 42 may be useful for other downstream tasks, e.g., computer vision or mobile manipulation [19], such  
 43 complex features may not be necessary for navigation tasks, e.g., when navigating human-occupied  
 44 spaces where everyone is respecting social etiquette [20]. Another argument for this limitation of  
 45 conventional computer vision models in real-world navigation is that humans only pay attention to  
 46 what is right in front of them to make decisions when navigating an environment. This efficiently  
 47 limits the observation region, as shown in Figure 1 right, which illustrates the difference between  
 48 features extracted by conventional computer vision models and those necessary to enable navigation  
 49 tasks.

50 Considering both the success of SSL on a variety of computer vision tasks and the oftentimes re-  
 51 dundant and confusing features provided by generic SSL models for navigation tasks, we present  
 52 a Vision-Action Navigation Pretraining (VANP) approach that completely relies on a pretext task  
 53 to train the visual encoder. The key observation behind VANP is when humans navigate crowded  
 54 spaces, we do not need to pay attention to all the people and objects in the scene, but only the ones  
 55 that affect our navigation trajectory *i.e.* observations that cause our actions. In this work, we re-  
 56 verse this causality by learning only relevant visual features with the help of future actions. To this  
 57 end, we leverage Barlow Twins’ redundancy-reduction principle [11] to train a visual encoder that  
 58 discards redundant features of an image for navigation using an action latent space (see Figure 2).  
 59 VANP focuses on extracting informative features from images that are aligned with the actions of a  
 60 navigating robot. Our experimental results suggest that VANP-extracted features are more informa-  
 61 tive for a downstream controller. We show that even a simple Convolutional Neural Network (CNN)  
 62 with 8 million parameters trained on only one million images can be as expressive for visual navi-  
 63 gation as a vision foundation model with 21 million parameters pre-trained on 142 million curated  
 64 images out of 1.2 billion source images.

65 The contributions of this work can be summarized as follows:

- 66 • We propose an SSL framework to train a visual encoder for robotic navigation tasks.
- 67 • We provide a concrete pre-trained SSL model for deployment in an end-to-end navigation
- 68 pipeline in social environments.

## 69 2 Related Work

70 Our work is motivated by several recent advances in Natural Language Processing (NLP) and com-  
71 puter vision, driven by the Self-Supervised Learning (SSL) paradigm. In this section, we categorize  
72 this pretraining paradigm into two groups for robotics and review related works pertaining to each.

73 **Pretraining for better representation:** General purpose models, also known as foundation models,  
74 pre-trained on pretext tasks can contribute to learning a rich representation that can help the model  
75 generalize to different downstream tasks in a zero/few-shot manner [16]. Foundation models for  
76 robot manipulation have been extensively studied in the literature [21, 22, 23, 24, 25, 26, 27]. For ex-  
77 ample, R3M [28] trained a general visual encoder for manipulation tasks on the Ego4D dataset [29],  
78 while CLIPort [30] leveraged the CLIP model [31] to enable language instructions for manipulation.  
79 Dadashi et al. [21] proposed AQuaDem, a framework to learn quantized actions from demonstra-  
80 tions in continuous action spaces, while VANP is doing the opposite by learning visual features from  
81 continuous action spaces. Luo et al. [22] improved AQuaDem by using VQ-VAE [32] for offline  
82 reinforcement learning. Huang et al. [23] proposed Skill Transformer to learn long-horizon robotic  
83 tasks with the help of transformers [33].

84 In autonomous driving Nazeri and Bohlouli [34] proposed two parallel networks, one encodes fea-  
85 tures from the past, and the other encodes plausible features from the future to expand the ob-  
86 servation window so the model can make well-informed decisions. Codevilla et al. [35] showed  
87 that a deeper model can play an important role in training better policies. It is apparent that  
88 most of the works in AVs use pre-trained computer vision models that are trained on ImageNet  
89 [36, 35, 37, 38, 34, 39, 40, 41, 42]. However, VANP shows that a visual encoder specific to naviga-  
90 tion tasks can help in learning better policies compared to pretraining with ImageNet.

91 **Pretraining for better policies:** Foundation models can not only help in learning a rich repre-  
92 sentation but also be used as a policy to generalize to multiple robotic tasks. SayCan [43] used  
93 Large Language Models (LLMs) to learn robotics skills by grounding LLMs and value functions  
94 in the physical world. Li et al. [44] and Reid et al. [45] used pre-trained LLMs as the policy back-  
95 bone. VPT [46] pseudo-labeled Minecraft YouTube videos to learn a behavior cloning policy that  
96 can craft diamonds. VPT learns the inverse dynamics while VANP uses dynamics to learn visual  
97 features. GNM [47] learned a general policy to drive any robot by combining multiple datasets  
98 of different robot types. ViNT [48] further improved GNM by replacing the policy network with  
99 a transformer [33]. To the best of our knowledge, no prior work has used actions as the pretext  
100 training signal to learn visual features for visual navigation.

101 **Large-Scale Datasets:** Large-scale datasets are the primary driver of recent advances in SSL. Data  
102 collection in computer vision is relatively straightforward compared to interaction-rich robotics nav-  
103 igation. One reason for this is that collecting a large-scale navigation dataset through human teleop-  
104 eration is expensive. Additionally, collecting interaction-rich datasets can be potentially dangerous  
105 due to the risk of collisions between humans and robots. Despite these challenges, the robotics com-  
106 munity has made tremendous efforts to collect interaction-rich datasets in the real world in recent  
107 years [49, 50, 51, 52]. SCAND [49] was one of the first efforts to collect navigation data in social  
108 environments at large by using teleoperated Spot and Jackal robots. MuSoHu [50] is another effort  
109 to collect 20 hours of human interactions in crowded spaces by using a human wearing a helmet  
110 equipped with different sensors. SANPO [52] also used humans to collect both real and synthetic  
111 datasets of nearly 15 hours of annotated videos for both vision tasks and robotics navigation. In  
112 this work, we combine both SCAND and MuSoHu of both robot and human navigation demonstra-  
113 tions respectively to create a dataset of nearly 1 million visual navigation samples with real-world  
114 human-robot interactions.

## 115 3 Methodology

116 This section formally defines the end-to-end visual navigation task and describes the Vision-Action  
117 Navigation Pretraining (VANP) procedure for the visual encoder.

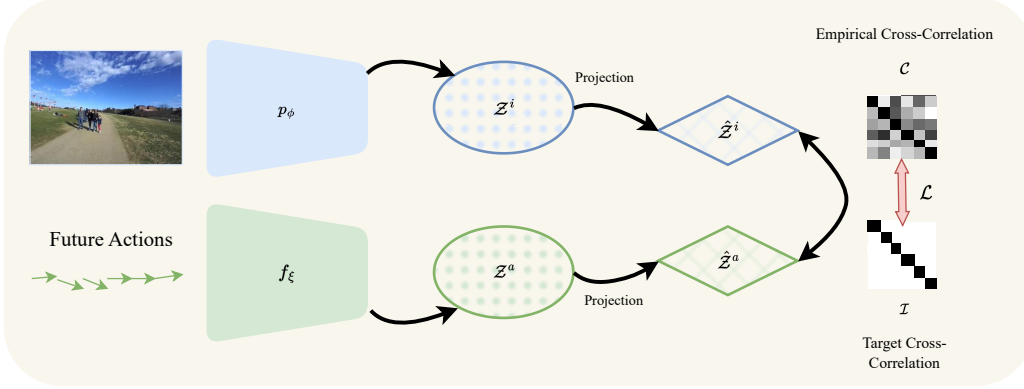


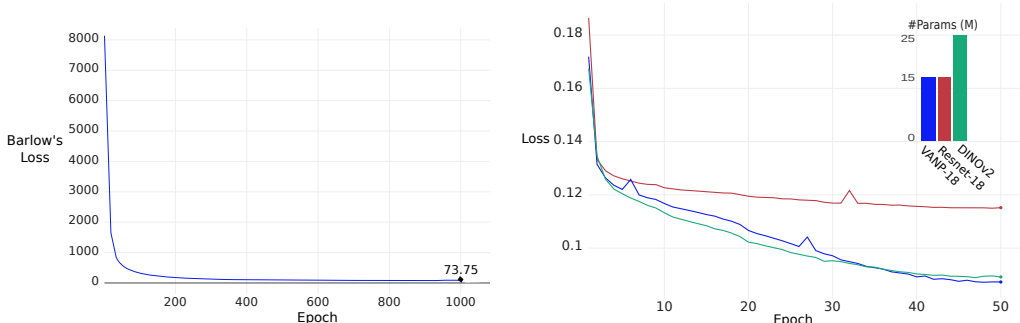
Figure 2: VANP first maps the sequence of actions to a sparser, higher-dimensional space  $Z^a$  (green). Then, leveraging Barlow Twins’ redundancy reduction principle, VANP induces sparsity on  $Z^i$  (blue), the visual embedding, by minimizing the mutual information between  $Z^i$  and  $Z^a$ .

118 **Problem Definition:** Visual navigation is the task of navigating an environment with only RGB  
 119 camera input. Unlike conventional geometric navigation tasks using, e.g., LiDARs or depth images,  
 120 this task is challenging due to the lack of explicit geometric information. The visual navigation  
 121 problem can be formalized as follows. **Input:** The robot is given a sequence of past and current  
 122 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  
 123 step,  $\tau_p$  is the number of past frames, and  $\mathcal{O}$  is the space of all possible image sequences. The robot  
 124 is also given its current goal e.g. GPS coordinates, pose, image, or next local coordinate in 2D space,  
 125  $g \in \mathcal{G}$ , which determines the direction it should move in the next time step. **Output:** The robot must  
 126 select an action,  $a_t \in \mathcal{A}$  consists of continuous linear and angular velocities, where  $\mathcal{A} \in [-1, 1]^2$  is  
 127 the action space, where  $[-1, 1]$  maps to the minimal and maximal linear and angular velocity of the  
 128 robot. **Visual Navigation:** The goal is to learn a policy,  $\pi_{\Theta} : \mathcal{O} \times \mathcal{G} \rightarrow \mathcal{A}$ , where  $\Theta$  represents the  
 129 policy’s parameters, to determine which action to take at each time step to reach its goal destination  
 130 efficiently while avoiding collisions with other agents and observing underlying social norms.

131 **End-To-End Model:** In end-to-end or holistic models we define the policy  $\pi_{\Theta}$  as follow:  $\mathbf{a} =$   
 132  $\pi_{\Theta}(\mathbf{o}, g) = \sigma_{\zeta}(p_{\phi}(\mathbf{o}) \oplus q_{\psi}(g))$ , where  $\sigma$  is the controller policy parameterized by  $\zeta$ ,  $p$  is the image  
 133 encoder parameterized by  $\phi$ ,  $q$  is the goal encoder parameterized by  $\psi$ , and  $\oplus$  is the concatenation of  
 134 two output vectors. To learn these parameters, two common approaches are (1) to learn all of them  
 135 together in an end-to-end manner which makes the training difficult and time-consuming or (2) to  
 136 train the image encoder separately and only fine-tune the goal encoder along with the controller to  
 137 reduce training time.

### 138 3.1 Vision-Action Model

139 VANP is inspired by the redundancy reduction principle of Barlow’s Twins to train the image en-  
 140 coder  $p$ . However, unlike vision self-supervised learning (SSL) models that work on the joint em-  
 141 bedding of augmented images [53, 54], VANP correlates the action space  $\mathcal{A}$  with the pixel latent  
 142 space  $\mathcal{O}$ . Under the assumption that every dynamic object or person in the environment will adhere  
 143 to social norms, we define VANP pretraining as follows: We sample a batch of  $(I^i, a_{t:t+\tau_F}^i)$  from  
 144 dataset  $D$  where  $i$  is the sample number,  $I^i$  is a single image at time  $t$  and  $a_{t:t+\tau_F}^i$  is a sequence of  
 145 actions starting from  $t$  and ending in  $t + \tau_F$ , where  $\tau_F$  is the number of frames in the future. We  
 146 then feed  $I^i$  to  $p_{\phi}$  and  $a_{t:t+\tau_F}^i$  to  $f_{\xi}$ , typically a multilayer perceptron (MLP), to learn image  $Z^i$  and  
 147 action  $Z^a$  embeddings, respectively. Finally, we use Barlow Twin’s objective function to learn  $\phi$   
 148 and  $\xi$ :



(a) Barlow’s loss smoothly decreases which leads to correlations between actions and visual features. (b) Downstream navigation training loss with Resnet-18, VANP-18, and DINOv2 as visual encoder.

Figure 3: Vision-Action Navigation Pretraining and Downstream Navigation Fine-Tuning.

$$\mathcal{L}_{BT} = \underbrace{\sum_i (1 - C_{ii})^2}_{\text{invariance term}} + \lambda \underbrace{\sum_i \sum_{j \neq i} C_{ij}^2}_{\text{redundancy reduction term}} \quad (1)$$

149 where  $\lambda$  is the trade-off between the first and second terms of the loss, and  $C$  is the cross-correlation  
 150 matrix computed between the outputs of the action and image embeddings along the batch dimen-  
 151 sion.

152 Leveraging Barlow’s objective function provides the advantage of not requiring negative samples,  
 153 which can be difficult to define in action space. For example, when a person is in front of us, there  
 154 may be two correct actions: overtake from the left or overtake from the right. Therefore, simply  
 155 negating the angular velocity does not give us a negative sample and may introduce ambiguity like  
 156 the example as mentioned earlier. Using actions from another sequence may not provide useful  
 157 information for visual navigation, as the actions are inherently conditioned on the observations.  
 158 Furthermore, the sparse high dimensional action latent space acts as a soft-whitening constraint on  
 159 the image latent space to reduce the redundancy in extracted features from the image [55, 11].

## 160 4 Preliminary Experimental Results

### 161 4.1 Implementation Details

162 We implement our method with PyTorch [56] and the training is performed on a single A100 GPU  
 163 with 80 gigabytes of memory. You can find the code here.

164 **Model architecture:** Considering the limited computation resources onboard most mobile robots,  
 165 we choose ResNet-18 [57] without the classification head as a lower latency image encoder and we  
 166 call it VANP-18. We use a multilayer perceptron (MLP) with two hidden layers as the action encoder  
 167 to produce the embeddings of  $Z^i, Z^a \in \mathbb{R}^{512}$ . Both encoders were followed by MLPs with three  
 168 layers as the projection heads to generate the final  $Z'^i, Z'^a \in \mathbb{R}^{8192}$ , the same as Zbontar et al. [11].  
 169 The most important challenge here is that the two distinct networks for producing the embeddings  
 170 have different modalities and therefore the output range significantly varies. Therefore, we initialize  
 171 all the deep networks with the Kaiming Normal initialization [57] with mean zero and variance one  
 172 to mitigate the discrepancy in the output of the networks. For the end-to-end model, we follow the  
 173 model used by Nguyen et al. [50]. We freeze the image encoder and only train the goal encoder and  
 174 controller during downstream task training for all the experiments.

175 **Optimization:** As proposed by Zbontar et al. [11], we use the LARS optimizer [58] and train the  
 176 model for 1000 epochs with a batch size of 16384. For the other hyperparameters, we use a learning  
 177 rate of 0.2 for the weights and 0.0048 for the biases. We use the first ten epochs as the warm-up  
 178 phase and update the learning rate by a factor of 8 during these epochs. We observe that, as suggested  
 179 by Zbontar et al. [11], using any other factor than 8 results in gradient explosion during training.



180 **Dataset:** We leverage two unique datasets: SCAND [49] and MuSoHu [50], both of which encap-  
 181 sulate robot and human navigation data from the egocentric perspective. Both large-scale real-world  
 182 datasets are collected in a variety of natural crowded public spaces. MuSoHu comprises approxi-  
 183 mately 20 hours of data captured from human egocentric motion. The recordings capture human  
 184 walking patterns in public spaces, providing insights for learning human-like, socially compliant  
 185 navigation behaviors. SCAND is an autonomous robot navigation dataset that captures 8.7 hours of  
 186 human-teleoperated robot navigation demonstrations in naturally crowded public spaces on a uni-  
 187 versity campus. By combining these two, we create a dataset of over 1 million samples to train the  
 188 image encoder on the pretext task. For pretext task training, we use a single image  $I_t \in \mathbb{R}^{70 \times 70}$   
 189 along with a sequence of actions  $a_{t:t+\tau_F} \in \mathbb{R}^{\tau_F \times 2}$  parsed at 25 Hz, comprising of 4 seconds in the  
 190 future. For the downstream task, we use a sequence of past observations  $I_{t-\tau_P:t} \in \mathbb{R}^{t \times 70 \times 70}$  along  
 191 with the polar coordinates of the next local goal  $g \in \mathbb{R}^2$  parsed at 4 Hz, containing 1.5 seconds  
 192 history as the network input to produce the actions  $\mathcal{A}_{t:t+\tau_F} \in \mathbb{R}^{\tau_F \times 2}$  for two seconds in the future.  
 193 In both stages, we use the augmentations proposed by Codevilla et al. [35].

## 194 4.2 Results Discussion

195 We report our preliminary results. To evaluate the effectiveness of VANP pretext training, we quan-  
 196 titatively compare it with DINOv2 [12], a self-supervised vision transformer model that learns uni-  
 197 versal features suitable for eight different visual tasks including depth estimation, semantic seg-  
 198 mentation, instance retrieval, dense and sparse matching. DINOv2 models exhibit robust out-of-  
 199 distribution performance, and the learned features can be used directly without fine-tuning. We use  
 200 DINOv2 as the upper performance bound and ResNet-18 trained on ImageNet ILSVRC-2012 as the  
 201 performance baseline. To ensure a fair comparison, the architectures of all other components of the  
 202 end-to-end model are kept fixed. During the downstream navigation task, we only train the goal en-  
 203 coder and controller, and the weights of the image encoder are frozen, regardless of the architecture  
 204 used.

205 We use ResNet-18 as the architecture for VANP pretext training (VANP-18). The smoothly de-  
 206 creasing Barlow’s loss in Figure 3a indicates that the learned visual features align with the actions.  
 207 Figure 3b shows the training loss for the downstream visual navigation task. During end-to-end  
 208 training, the weights of the visual encoder  $\phi$  are frozen, and we only train the goal encoder  $q_\psi$  and  
 209 the controller  $\sigma_\zeta$ . As shown in the figure, VANP-18 outperforms ResNet-18 with the same archi-  
 210 tecture but a different training paradigm which shows the effectiveness of VANP’s pretext training.  
 211 VANP-18 also achieves slightly better performance than DINOv2, a vision transformer trained on  
 212 billions of images. Figure 3b shows a comparison of parameters for the holistic model when we use  
 213 different visual encoders.

214 We separate a few trajectories from the dataset and use them as unseen scenarios for qualitative  
 215 evaluation. After qualitatively comparing the model outputs, we observe that VANP-18 performs  
 216 human avoidance. Figure 4 shows a few examples from the evaluation set. Note that negative values  
 217 for angular velocity denote turning right, and positive values mean turning left. VANP-18’s decisions  
 218 do not agree with human decisions in some scenarios, but it does not mean that they are completely  
 219 wrong: for example, the disagreement in the second red border image is because the human predicts  
 220 the group’s future trajectory, and decides to move from the left while VANP-18 decides to go from  
 221 the right.

## 222 5 Conclusions and Future Work

223 In this work, we propose a self-supervised training approach to train visual encoder models specifi-  
 224 cally designed for visual navigation. This approach is motivated by the observation that humans  
 225 only pay attention to a small region of their frontal view to make navigation decisions. By reversing  
 226 this observation, we use the decisions to extract only visual features that are relevant to the visual  
 227 navigation task, unlike computer vision models that tend to extract every detail in the environment,  
 228 which can lead to confusion of neural-based controllers.



Figure 4: VANP’s qualitative performance on unseen scenarios. Green: VANP outputs align with the demonstrations; Red: VANP outputs do not align with the demonstrations.

229 We hope that this work will inspire new ideas in the field of visual navigation. In the future, we  
 230 plan to conduct more real-world experiments and train deeper models, such as VAM-34 and VAM-  
 231 50, based on ResNet-34 and ResNet-50, respectively. In this work, we only use datasets collected  
 232 in social environments. Another future direction is to merge datasets from different environments,  
 233 such as off-road, indoor, outdoor, and social environments, to evaluate the generalizability of the  
 234 proposed approach.

235 **References**

- 236 [1] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. De-  
237 hghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby. An Image is Worth  
238 16x16 Words: Transformers for Image Recognition at Scale.
- 239 [2] Z. Gan, L. Li, C. Li, L. Wang, Z. Liu, and J. Gao. Vision-Language Pre-training: Basics,  
240 Recent Advances, and Future Trends.
- 241 [3] H. Cao, C. Tan, Z. Gao, G. Chen, P.-A. Heng, and S. Z. Li. A Survey on Generative Diffusion  
242 Model.
- 243 [4] W. Zhang, L. He, H. Wang, L. Yuan, and W. Xiao. Multiple Self-Supervised Auxiliary Tasks  
244 for Target-Driven Visual Navigation Using Deep Reinforcement Learning. 25(7):1007. ISSN  
245 1099-4300. doi:10.3390/e25071007.
- 246 [5] A. Lacoste, A. Luccioni, V. Schmidt, and T. Dandres. Quantifying the Carbon Emissions of  
247 Machine Learning.
- 248 [6] E. Strubell, A. Ganesh, and A. McCallum. Energy and Policy Considerations for Deep Learn-  
249 ing in NLP.
- 250 [7] D. Patterson, J. Gonzalez, Q. Le, C. Liang, L.-M. Munguia, D. Rothchild, D. So, M. Texier,  
251 and J. Dean. Carbon Emissions and Large Neural Network Training. . doi:10.48550/ARXIV.  
252 2104.10350.
- 253 [8] D. Patterson, J. Gonzalez, U. Hölzle, Q. H. Le, C. Liang, L.-M. Munguia, D. Rothchild, D. So,  
254 M. Texier, and J. Dean. The Carbon Footprint of Machine Learning Training Will Plateau,  
255 Then Shrink, .
- 256 [9] L. B. Heguerte, A. Bugeau, and L. Lanelongue. How to estimate carbon footprint when  
257 training deep learning models? A guide and review.
- 258 [10] T. Chen, S. Kornblith, K. Swersky, M. Norouzi, and G. Hinton. Big self-supervised models are  
259 strong semi-supervised learners.
- 260 [11] J. Zbontar, L. Jing, I. Misra, Y. LeCun, and S. Deny. Barlow Twins: Self-Supervised Learning  
261 via Redundancy Reduction.
- 262 [12] M. Oquab, T. Darcet, T. Moutakanni, H. Vo, M. Szafraniec, V. Khalidov, P. Fernandez, D. Haz-  
263 iza, F. Massa, A. El-Nouby, M. Assran, N. Ballas, W. Galuba, R. Howes, P.-Y. Huang, S.-W. Li,  
264 I. Misra, M. Rabbat, V. Sharma, G. Synnaeve, H. Xu, H. Jegou, J. Mairal, P. Labatut, A. Joulin,  
265 and P. Bojanowski. DINOv2: Learning Robust Visual Features without Supervision.
- 266 [13] A. Bardes, J. Ponce, and Y. LeCun. MC-JEPA: A Joint-Embedding Predictive Architecture for  
267 Self-Supervised Learning of Motion and Content Features.
- 268 [14] S. Gidaris, P. Singh, and N. Komodakis. Unsupervised representation learning by predicting  
269 image rotations.
- 270 [15] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick. Masked autoencoders are scalable  
271 vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern  
272 Recognition*, pages 16000–16009.
- 273 [16] R. Balestriero, M. Ibrahim, V. Sobal, A. Morcos, S. Shekhar, T. Goldstein, F. Bordes,  
274 A. Bardes, G. Mialon, Y. Tian, A. Schwarzschild, A. G. Wilson, J. Geiping, Q. Garrido, P. Fer-  
275 nandez, A. Bar, H. Pirsiavash, Y. LeCun, and M. Goldblum. A Cookbook of Self-Supervised  
276 Learning.



- 277 [17] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A large-scale hierar-  
278 chical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*,  
279 pages 248–255. doi:10.1109/CVPR.2009.5206848.
- 280 [18] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zit-  
281 nick. Microsoft COCO: Common Objects in Context. In D. Fleet, T. Pajdla, B. Schiele,  
282 and T. Tuytelaars, editors, *Computer Vision – ECCV 2014*, volume 8693, pages 740–  
283 755. Springer International Publishing. ISBN 978-3-319-10601-4 978-3-319-10602-1. doi:  
284 10.1007/978-3-319-10602-1\_48.
- 285 [19] S. Karamcheti, S. Nair, A. S. Chen, T. Kollar, C. Finn, D. Sadigh, and P. Liang. Language-  
286 Driven Representation Learning for Robotics.
- 287 [20] A. Robicquet, A. Sadeghian, A. Alahi, and S. Savarese. Learning Social Etiquette: Human  
288 Trajectory Understanding In Crowded Scenes. In B. Leibe, J. Matas, N. Sebe, and M. Welling,  
289 editors, *Computer Vision – ECCV 2016*, volume 9912, pages 549–565. Springer International  
290 Publishing. ISBN 978-3-319-46483-1 978-3-319-46484-8. doi:10.1007/978-3-319-46484-8-  
291 33.
- 292 [21] R. Dadashi, L. Hussenot, D. Vincent, S. Girgin, A. Raichuk, M. Geist, and O. Pietquin. Con-  
293 tinuous control with action quantization from demonstrations. In *International Conference on*  
294 *Machine Learning*, pages 4537–4557. PMLR.
- 295 [22] J. Luo, P. Dong, J. Wu, A. Kumar, X. Geng, and S. Levine. Action-quantized offline reinforce-  
296 ment learning for robotic skill learning. In *7th Annual Conference on Robot Learning*.
- 297 [23] X. Huang, D. Batra, A. Rai, and A. Szot. Skill Transformer: A Monolithic Policy for Mobile  
298 Manipulation.
- 299 [24] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan, K. Haus-  
300 man, A. Herzog, J. Hsu, J. Ibarz, B. Ichter, A. Irpan, T. Jackson, S. Jesmonth, N. J. Joshi,  
301 R. Julian, D. Kalashnikov, Y. Kuang, I. Leal, K.-H. Lee, S. Levine, Y. Lu, U. Malla, D. Manju-  
302 nath, I. Mordatch, O. Nachum, C. Parada, J. Peralta, E. Perez, K. Pertsch, J. Quiambao, K. Rao,  
303 M. Ryoo, G. Salazar, P. Sanketi, K. Sayed, J. Singh, S. Sontakke, A. Stone, C. Tan, H. Tran,  
304 V. Vanhoucke, S. Vega, Q. Vuong, F. Xia, T. Xiao, P. Xu, S. Xu, T. Yu, and B. Zitkovich. RT-1:  
305 Robotics Transformer for Real-World Control at Scale.
- 306 [25] N. Di Palo, A. Byravan, L. Hasenclever, M. Wulfmeier, N. Heess, and M. Riedmiller. Towards  
307 A Unified Agent with Foundation Models. doi:10.48550/ARXIV.2307.09668.
- 308 [26] J. Carvalho, A. T. Le, M. Baierl, D. Koert, and J. Peters. Motion Planning Diffusion: Learning  
309 and Planning of Robot Motions with Diffusion Models.
- 310 [27] A. Hiranaka, M. Hwang, S. Lee, C. Wang, L. Fei-Fei, J. Wu, and R. Zhang. Primitive Skill-  
311 based Robot Learning from Human Evaluative Feedback.
- 312 [28] S. Nair, A. Rajeswaran, V. Kumar, C. Finn, and A. Gupta. R3M: A universal visual represen-  
313 tation for robot manipulation. In *Conference on Robot Learning*, pages 892–909. PMLR.
- 314 [29] K. Grauman, A. Westbury, E. Byrne, Z. Chavis, A. Furnari, R. Girdhar, J. Hamburger, H. Jiang,  
315 M. Liu, X. Liu, M. Martin, T. Nagarajan, I. Radosavovic, S. K. Ramakrishnan, F. Ryan,  
316 J. Sharma, M. Wray, M. Xu, E. Z. Xu, C. Zhao, S. Bansal, D. Batra, V. Cartillier, S. Crane,  
317 T. Do, M. Doulaty, A. Erapalli, C. Feichtenhofer, A. Fragomeni, Q. Fu, C. Fuegen, A. Ge-  
318 breselasie, C. Gonzalez, J. Hillis, X. Huang, Y. Huang, W. Jia, W. Khoo, J. Kolar, S. Kottur,  
319 A. Kumar, F. Landini, C. Li, Y. Li, Z. Li, K. Mangalam, R. Modhugu, J. Munro, T. Murrell,  
320 T. Nishiyasu, W. Price, P. R. Puentes, M. Ramazanova, L. Sari, K. Somasundaram, A. Souther-  
321 land, Y. Sugano, R. Tao, M. Vo, Y. Wang, X. Wu, T. Yagi, Y. Zhu, P. Arbelaez, D. Crandall,  
322 D. Damen, G. M. Farinella, B. Ghanem, V. K. Ithapu, C. V. Jawahar, H. Joo, K. Kitani, H. Li,

- 323 R. Newcombe, A. Oliva, H. S. Park, J. M. Rehg, Y. Sato, J. Shi, M. Z. Shou, A. Torralba,  
324 L. Torresani, M. Yan, and J. Malik. Ego4D: Around the World in 3,000 Hours of Egocentric  
325 Video. In *IEEE/CVF Computer Vision and Pattern Recognition (CVPR)*.
- 326 [30] M. Shridhar, L. Manuelli, and D. Fox. Cliport: What and where pathways for robotic manipu-  
327 lation. In *Conference on Robot Learning*, pages 894–906. PMLR.
- 328 [31] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell,  
329 P. Mishkin, J. Clark, G. Krueger, and I. Sutskever. Learning transferable visual models from  
330 natural language supervision. In M. Meila and T. Zhang, editors, *Proceedings of the 38th Inter-  
331 national Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning  
332 Research*, pages 8748–8763. PMLR.
- 333 [32] A. van den Oord, O. Vinyals, and K. Kavukcuoglu. Neural discrete representation learning.
- 334 [33] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polo-  
335 sukhin. Attention is all you need. 30.
- 336 [34] M. H. Nazeri and M. Bohlouli. Exploring Reflective Limitation of Behavior Cloning in Au-  
337 tonomous Vehicles. In *2021 IEEE International Conference on Data Mining (ICDM)*, pages  
338 1252–1257. IEEE. ISBN 978-1-66542-398-4. doi:10.1109/ICDM51629.2021.00153.
- 339 [35] F. Codevilla, E. Santana, A. M. Lopez, and A. Gaidon. Exploring the Limitations of Behavior  
340 Cloning for Autonomous Driving. In *The IEEE International Conference on Computer Vision  
341 (ICCV)*. IEEE. doi:10.1109/ICCV.2019.00942.
- 342 [36] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun. CARLA: An open urban driving  
343 simulator. In *Proceedings of the 1st Annual Conference on Robot Learning*, pages 1–16.
- 344 [37] D. Chen, B. Zhou, V. Koltun, and P. Krähenbühl. Learning by cheating. In L. P. Kaelbling,  
345 D. Kragic, and K. Sugiura, editors, *Proceedings of the Conference on Robot Learning*, volume  
346 100 of *Proceedings of Machine Learning Research*, pages 66–75. PMLR, 2020-11-30, 2020.
- 347 [38] E. Ohn-Bar, A. Prakash, A. Behl, K. Chitta, and A. Geiger. Learning Situational Driving.  
348 In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE.  
349 doi:10.1109/cvpr42600.2020.01131.
- 350 [39] D. Chen, V. Koltun, and P. Krähenbühl. Learning to drive from a world on rails. In *ICCV*.
- 351 [40] K. Chitta, A. Prakash, B. Jaeger, Z. Yu, K. Renz, and A. Geiger. TransFuser: Imitation with  
352 transformer-based sensor fusion for autonomous driving.
- 353 [41] A. Hu, G. Corrado, N. Griffiths, Z. Murez, C. Gurau, H. Yeo, A. Kendall, R. Cipolla, and  
354 J. Shotton. Model-Based Imitation Learning for Urban Driving.
- 355 [42] B. Jaeger, K. Chitta, and A. Geiger. Hidden biases of end-to-end driving models. In *Interna-  
356 tional Conference on Computer Vision (ICCV)*.
- 357 [43] M. Ahn, A. Brohan, N. Brown, Y. Chebotar, O. Cortes, B. David, C. Finn, C. Fu, K. Gopalakr-  
358 ishnan, K. Hausman, et al. Do as i can, not as i say: Grounding language in robotic affordances.
- 359 [44] S. Li, X. Puig, C. Paxton, Y. Du, C. Wang, L. Fan, T. Chen, D.-A. Huang, E. Akyürek,  
360 A. Anandkumar, et al. Pre-trained language models for interactive decision-making. 35:  
361 31199–31212.
- 362 [45] M. Reid, Y. Yamada, and S. S. Gu. Can wikipedia help offline reinforcement learning?
- 363 [46] B. Baker, I. Akkaya, P. Zhokov, J. Huizinga, J. Tang, A. Ecoffet, B. Houghton, R. Sampedro,  
364 and J. Clune. Video pretraining (vpt): Learning to act by watching unlabeled online videos.  
365 35:24639–24654.

- 366 [47] D. Shah, A. Sridhar, A. Bhorkar, N. Hirose, and S. Levine. GNM: A General Navigation Model  
367 to Drive Any Robot, .
- 368 [48] D. Shah, A. Sridhar, N. Dashora, K. Stachowicz, K. Black, N. Hirose, and S. Levine. ViNT: A  
369 Foundation Model for Visual Navigation, .
- 370 [49] H. Karnan, A. Nair, X. Xiao, G. Warnell, S. Pirk, A. Toshev, J. Hart, J. Biswas, and P. Stone.  
371 Socially Compliant Navigation Dataset (SCAND): A Large-Scale Dataset of Demonstrations  
372 for Social Navigation.
- 373 [50] D. M. Nguyen, M. Nazeri, A. Payandeh, A. Datar, and X. Xiao. Toward human-like social  
374 robot navigation: A large-scale, multi-modal, social human navigation dataset.
- 375 [51] N. Hirose, D. Shah, A. Sridhar, and S. Levine. SACSoN: Scalable Autonomous Data Collec-  
376 tion for Social Navigation.
- 377 [52] S. M. Waghmare, K. Wilber, D. Hawkey, X. Yang, M. Wilson, S. Debats, C. Nuengsigkapijan,  
378 A. Sharma, L. Pandikow, H. Wang, H. Adam, and M. Sirotenko. SANPO: A Scene Under-  
379 standing, Accessibility, Navigation, Pathfinding, Obstacle Avoidance Dataset.
- 380 [53] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton. A simple framework for contrastive learning  
381 of visual representations. In *International Conference on Machine Learning*, pages 1597–  
382 1607. PMLR.
- 383 [54] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick. Momentum contrast for unsupervised visual  
384 representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and  
385 Pattern Recognition*, pages 9729–9738.
- 386 [55] A. Ermolov, A. Siarohin, E. Sangineto, and N. Sebe. Whitening for self-supervised repre-  
387 sentation learning. In M. Meila and T. Zhang, editors, *Proceedings of the 38th International  
388 Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*,  
389 pages 3015–3024. PMLR.
- 390 [56] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin,  
391 N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Te-  
392 jani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala. PyTorch: An imperative  
393 style, high-performance deep learning library. In *Advances in Neural Information Processing  
394 Systems 32*, pages 8024–8035. Curran Associates, Inc.
- 395 [57] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. pages  
396 770–778. doi:10.1109/cvpr.2016.90.
- 397 [58] Y. You, I. Gitman, and B. Ginsburg. Large batch training of convolutional networks. *arXiv  
398 preprint arXiv:1708.03888*, 2017.