

Robust Manipulation with Spatial Features

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** Our goal is to develop visual pre-training strategies that enable more
2 robust and efficient manipulation policy learning. We find that a Vision Trans-
3 former trained with a distillation loss that biases representations towards shape
4 exhibits strong zero-shot transfer performance on the kitchen shift suite, even
5 when compared to baselines trained on larger and more task-relevant datasets.
6 When finetuned, the attention heads of a transformer trained with a shape bias can
7 be visualized as a spatial feature map, which emergently segments manipulation-
8 relevant objects in an image. By leveraging each of these insights, we are able to
9 improve the average zero-shot performance of policies trained on the sliding door
10 task within the FrankaKitchen environment by nearly 2x compared to the next best
11 method. Additionally, we are able to improve maximum success in distribution by
12 13% by masking out attention heads that attend to distractors.

13 **Keywords:** Manipulation, visual pre-training, self-supervision

14 1 Introduction

15 Our goal is to learn robotic manipulation policies from images. For many computer vision tasks,
16 models can be applied off-the-shelf in new environments with little to no task-specific tuning. In
17 spite of this success, robotic policies learned from pixels remain surprisingly brittle. One common
18 approach to learning from pixels follows a formula that is familiar to many computer vision practi-
19 tioners: pre-train a self-supervised network on a broad and diverse image dataset before fine-tuning
20 on task specific data and labels. We expect this strategy to yield a model that’s capable of general-
21 izing the downstream task across visually diverse environments, but when roboticists try this same
22 formula, learned policies break in the presence of a distractor, under subtle lighting changes, and
23 after a slight change in the camera position.

24 Recent work posits that the missing piece is a large dataset of object interactions across diverse
25 environments—the ImageNet or CommonCrawl of manipulation. Indeed, training on large datasets
26 of first person human interaction data increases policy performance downstream. However, these
27 policies remain brittle to even small distribution shifts that commonly occur during deployment.
28 Why are robotic policies learned from visual features so sensitive to distribution shift compared to
29 other tasks that rely on visual information?

30 Successful manipulation requires spatial reasoning. To that end, past work introduced structured
31 representations (e.g., keypoints) that capture the spatial aspects of the visual observation space at the
32 cost of expressivity. Instead of introducing explicit structure into the representation, we leverage an
33 encoder pre-trained with a self-supervised distillation loss (DiNo) [1] that biases the representation
34 towards shape. We show that policies learned on top of these representations are more robust in
35 the presence of visual distribution shift even when compared to representations learned from larger
36 and more task-relevant datasets. Unlike more structured approaches, pre-training with DiNo can be
37 applied to any architecture or dataset and doesn’t require explicit supervision.

38 Another benefit of this encoder choice is that the transformer attention heads can be visualized as
39 a spatial heatmap. We find that the visualized attention heads can be interacted with in predictable

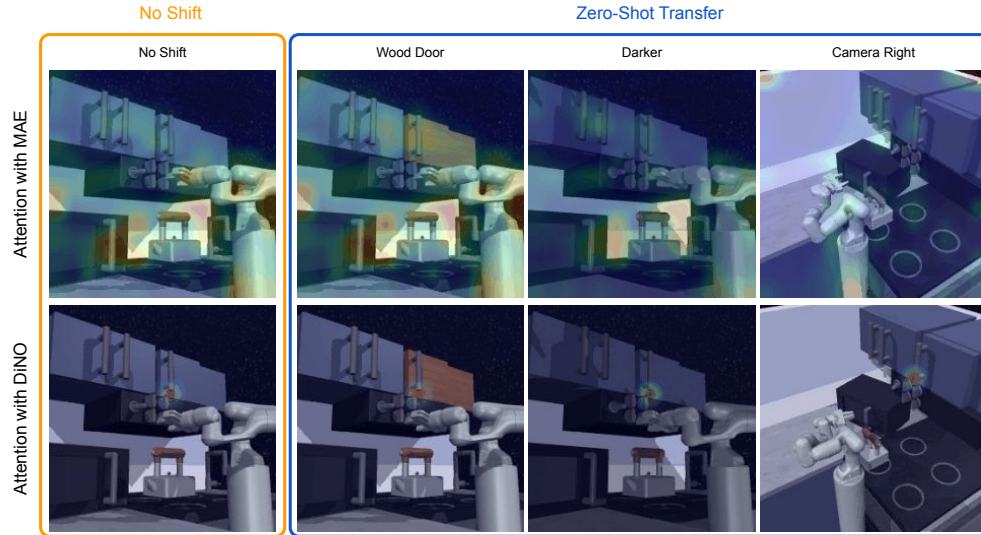


Figure 1: Encoders trained with losses that induce spatial heatmaps transfer to new viewpoints and textures zero-shot.

40 ways. Concretely, by masking out attention heads that segment task irrelevant parts of the image,
 41 we can improve the average performance of a policy trained on top of this architecture. This hints at
 42 the possibility of using the segmentation performance of the attention heads as a heuristic to identify
 43 good models to serve as representations for manipulation.

44 This work demonstrates that pre-training a visual encoder with a loss that is biased towards shape
 45 can dramatically improve policy performance under distribution shifts. We demonstrate this across
 46 27 different training settings of the door sliding task within the FrankaKitchen [2] environment and
 47 evaluate across 7 different texture and lighting changes adapted from the Kitchen Shift benchmark
 48 [3]. We find that a vision transformer trained with a shape-biased distillation loss (DiNo) strongly
 49 outperforms both Transformer [4] and ResNet [5] architectures trained on much larger and task
 50 relevant datasets. The attention heads from a shape-biased Vision Transformer can be visualized as
 51 spatial heatmaps, which can be visually inspected to identify task relevant and irrelevant heads. We
 52 explore this idea by masking attention heads that attend to irrelevant background pixels and observe
 53 a 13% boost to policy performance. After fine-tuning on a robotics task, these attention maps also
 54 emergently highlight task-relevant keypoints in the scene that are robust across visual distribution
 55 shifts. This is true for even dramatic visual shifts such as camera viewpoint. We believe that shape-
 56 biased losses could be a new standard for pre-training visual encoders for manipulation and that
 57 leveraging the attention heads of Vision Transformers (e.g., by masking) could lead to further policy
 58 improvements.

59 2 Related Work

60 **Policy adaptation.** Policies learned from pixels are known to be sensitive to distractors. Policy
 61 adaptation approaches aim to resolve this instability by continuing to train self-supervised visual
 62 representations between during deployment [6], improving the transferability of encoders through
 63 augmentations [7, 8], or collecting exploration data in the target environment to align source and
 64 target representations [9]. Unlike this work, our method doesn't require any target domain data or
 65 hand-designed augmentations.

66 **Representation learning for manipulation.** The correct approach to visual representation learn-
 67 ing for robotics is still an open question. Some work has analysed the transfer quality of a vari-
 68 ety different supervised vision tasks to robotics tasks [10]. Unlike this work, training with DiNo

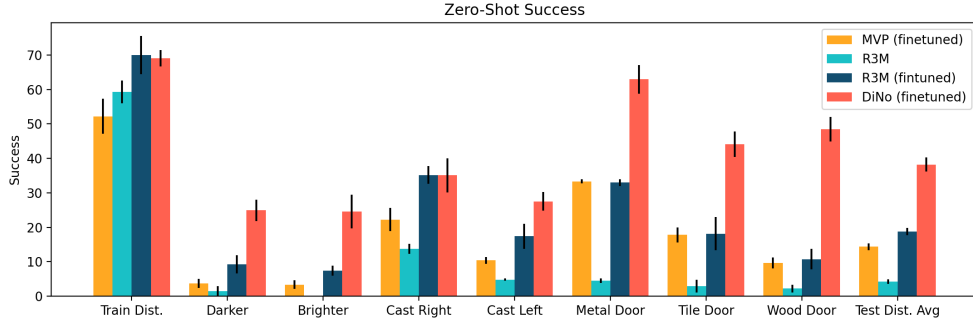


Figure 2: Zero shot performance on the kitchen shift suite.

69 does not require any labels and so it can be readily adapted to more robotics-relevant datasets. We
 70 compare directly against works that have developed self-supervised losses for manipulation [11]
 71 on manipulation-relevant datasets [12] or directly evaluated existing self-supervision approaches on
 72 such datasets [13].

73 3 Robust Manipulation with Spatial Features

74 We study two questions: (1) Can encoders trained with a shape-biased loss perform better under
 75 visual distribution shifts than other self-supervised losses? (2) Can the intuitive interpretations of
 76 attention map visualizations be leveraged to improve policy performance during training?

77 **Shape bias improves zero-shot transfer.** We follow the same evaluation protocol as R3M on
 78 the sliding door task. On top of each encoder, we train a two-layer MLP with imitation learning
 79 to perform the sliding door task. We compare across 3 seeds, 3 levels of demonstrations, and 3
 80 camera angles. We then evaluate the performance of the policy and encoder across a subset of
 81 visual distribution shifts in the Kitchen Shift benchmark. This includes changing lighting—making
 82 the lighting darker, making the lighting brighter, lighting cast left, and lighting cast right—as well
 83 as changing the texture wrapping the cabinet of the sliding door to be wood, metal, or tile. The
 84 R3M training and testing environments modify FrankaKitchen by randomizing the position of the
 85 kitchen, so we reimplement these distribution shifts in the R3M evaluation environment. Because
 86 the kitchen position is randomized, the task is much more difficult to solve using memorization. We
 87 expect replay data to perform much worse than in the original Kitchen Shift benchmark.

88 We compare a Vision Transformer trained with a shape-biased loss (DiNo) against three other visual
 89 representation learning approaches. In MVP we borrow the encoder from Radosavovic et al. [13]
 90 and finetune. MVP leverages a ViT trained with masked autoencoding (MAE) on a mixture of
 91 human interaction data including Ego-4D. We also compare against a frozen and finetuned model
 92 from R3M [11]. R3M utilizes a ResNet-50 architecture trained on top of Ego-4D.

93 The zero-shot performance of each model across distribution shifts can be found in Figure 2. On
 94 the left, we present the performance of each model without any distribution shift. We then plot the
 95 performance the models by shift type and show the performance averaged across shifts on the right.
 96 For all of the shifts, we average results across level of demonstrations and camera angles and then
 97 take the average and standard error over seeds.

98 Visualizations of the attention heads after training are presented in Figure 1. Similar to DiNo, we
 99 visualize attention heads by mapping the weight at each head at the output of the last block to a
 100 heat map and smoothing the final map with bilinear interpolation. The finetuned MVP model is
 101 visualized on the top and the fine-tuned DiNo model is visualized on the bottom. Of the six heads
 102 in the last block, we select the head that best attends to the manipulation relevant objects by visual
 103 inspection. At the best head, fine-tuning with DiNo appears to give more manipulation-relevant and

104 precise attention heads. Surprisingly, the same attention head is consistent across texture shifts and
105 across viewpoints.

106 **Leveraging attention for better policy learning.** In this section, we study the question: are the
107 attention head visualizations useful to the extent that they enable the development of a better per-
108 forming policy? This is an important proof-of-concept that opens the door for future work to improve
109 policy performance by leveraging attention heads that segment objects that are relevant to the de-
110 sired manipulation task. For example, a practitioner could decide to mask a head that attends to the
111 microwave if the policy needs to open the cabinet.

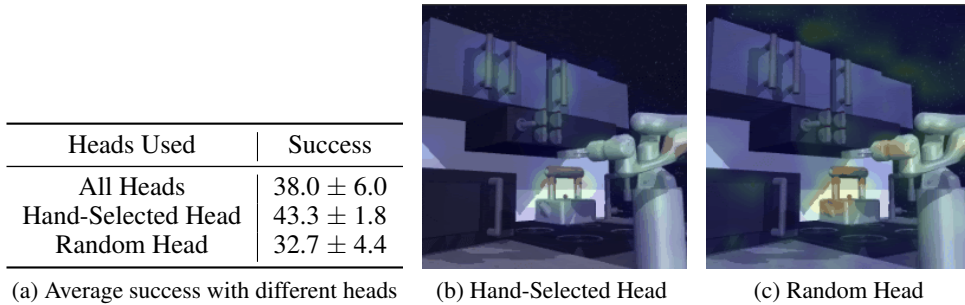


Figure 3: (Left) If we mask out all but an intuitively-correct hand-selected head, we can boost average policy performance by 13%. (Middle) A visualization of an attention head before fine-tuning on the target task. Without any environment data, the attention head segments manipulation relevant objects. (Right) A visualization of an attention head selected at random. Compared to the hand selected head, the random head segments the background, which is irrelevant to the sliding door task.

112 We focus on the sliding door task trained with 5 human demonstrations and the left camera view-
113 point. We visually inspect each of the 6 attention heads of a DiNo-pretrained vision transformer
114 and select the head that segments the most task-relevant objects. We mask all but the hand-selected
115 head and compare the success of training an MLP without finetuning after 1000 training steps. We
116 present the average performance results with standard error across 3 seeds in Table 3a. For an addi-
117 tional baseline, we also report the results of masking all but a random head. Attention heads are
118 masked by zeroing out the weights that map from input vectors to query, key, and value vectors. We
119 only visualize and mask heads at the last attention block. After masking out the hand-selected head,
120 success after 1000 training steps sees a modest performance improvement with reduced variance
121 compared to using all heads.

122 4 Conclusion

123 In this paper we studied two questions related to visual representation learning for manipulation.
124 First, we find that pre-training with a loss that induces a shape bias can provide strong performance
125 gains when evaluating policies under visual distribution shift. Second, we present a proof of concept
126 that leverages the insight that the attention heads of a DiNo-trained Vision Transformer segment task
127 relevant objects. Our findings open up important questions for future work, such as: could training
128 larger and more task-relevant datasets, such as Ego-4D, with a shape-biased loss further improve
129 policy learning performance?

References

- [1] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2021.
- [2] A. Gupta, V. Kumar, C. Lynch, S. Levine, and K. Hausman. Relay policy learning: Solving long-horizon tasks via imitation and reinforcement learning. In *Conference on Robot Learning*, pages 1025–1037. PMLR, 2020.
- [3] E. Xing, A. Gupta, S. Powers*, and V. Dean*. Kitchenshift: Evaluating zero-shot generalization of imitation-based policy learning under domain shifts. In *NeurIPS 2021 Workshop on Distribution Shifts: Connecting Methods and Applications*, 2021. URL <https://openreview.net/forum?id=DdglKo8hBq0>.
- [4] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need. 2017. URL <https://arxiv.org/pdf/1706.03762.pdf>.
- [5] 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)*, pages 770–778, 2016.
- [6] N. Hansen, R. Jangir, Y. Sun, G. Alenyà, P. Abbeel, A. A. Efros, L. Pinto, and X. Wang. Self-supervised policy adaptation during deployment. In *International Conference on Learning Representations*, 2021.
- [7] N. Hansen and X. Wang. Generalization in reinforcement learning by soft data augmentation. In *International Conference on Robotics and Automation*, 2021.
- [8] L. Fan, G. Wang, D.-A. Huang, Z. Yu, L. Fei-Fei, Y. Zhu, and A. Anandkumar. Secant: Self-expert cloning for zero-shot generalization of visual policies. In M. Meila and T. Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 3088–3099. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/fan21c.html>.
- [9] T. Yoneda, G. Yang, M. R. Walter, and B. Stadie. Invariance through latent alignment, 2021.
- [10]
- [11] S. Naie, A. Rajeswaran, V. Kumar, C. Finn, and A. Gupta. R3m: A universal visual representation for robot manipulation. In *Conference on Robot Learning (CoRL)*, 2022.
- [12] K. Grauman, A. Westbury, E. Byrne, Z. Q. 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. González, J. M. Hillis, X. Huang, Y. Huang, W. Jia, W. Y. H. Khoo, J. Kolár, 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. Ramazanov, L. Sari, K. K. Somasundaram, A. Southerland, Y. Sugano, R. Tao, M. Vo, Y. Wang, X. Wu, T. Yagi, Y. Zhu, P. Arbeláez, D. J. Crandall, D. Damen, G. M. Farinella, B. Ghanem, V. K. Ithapu, C. V. Jawahar, H. Joo, K. Kitani, H. Li, R. A. 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. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18973–18990, 2022.
- [13] I. Radosavovic, T. Xiao, S. James, P. Abbeel, J. Malik, and T. Darrell. Real-world robot learning with masked visual pre-training. *CoRL*, 2022.