

On Pre-Training for Visuo-Motor Control: Revisiting a Learning-from-Scratch Baseline

Anonymous Author(s)
Affiliation
Address
email

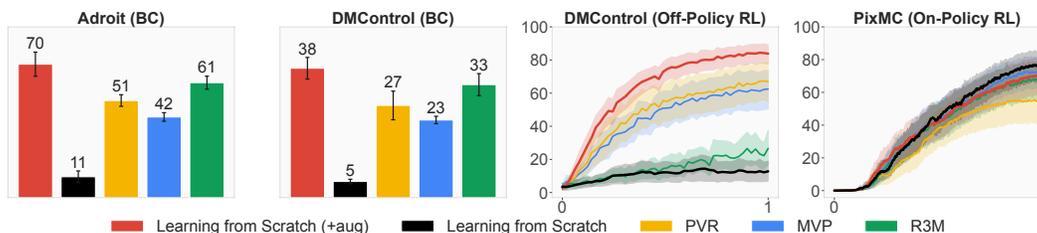


Figure 1: **Pre-training vs. Learning-from-Scratch (LfS)**. Success rate (Adroit, PixMC) and normalized return (DMControl) in each of the three task domains that we consider (aggregated across tasks). BC results are averages of top-3 evaluations over 100 epochs [1], and RL results are reported as a function of environment steps [2, 3], normalized to $[0, 1]$ since number of steps differ between tasks. We evaluate strong LfS baselines [2, 4] and find them to be competitive with recent frozen pre-trained representations. We report mean and 95% confidence intervals over 5 seeds.

1 **Abstract:** We revisit a simple Learning-from-Scratch baseline for visuo-motor
2 control that uses data augmentation and a shallow ConvNet. We find that this
3 baseline has competitive performance with recent methods that leverage frozen
4 visual representations trained on large-scale vision datasets.

5 **Keywords:** to pre-train; not to pre-train

6 1 Introduction

7 Large-scale pre-training has delivered promising results in computer vision [5, 6, 7, 8] and natural
8 language processing [9, 10, 11, 12]. Recent works have extended the pre-training paradigm to visuo-
9 motor control by leveraging pre-trained visual representations for policy learning [1, 13, 3, 14, 15].
10 These works train a visual representation using large out-of-domain vision datasets like ImageNet [16]
11 and Ego4D [17], and freeze the vision model weights for downstream policy learning. When
12 compared to simple Learning-from-Scratch (LfS) methods for visuo-motor control, these works
13 find that frozen pre-trained representations help achieve high sample efficiency and/or asymptotic
14 performance across a variety of domains and algorithms.

15 However, there exists a rich body of work on strategies to improve performance of LfS methods,
16 such as auxiliary self-supervised representation learning [18, 19] or using carefully curated data
17 augmentations [20, 2, 4]. To gain a sharp understanding of the advantages of visual pre-training
18 for motor control, it is necessary to establish strong LfS baselines. Towards this end, we take the
19 experimental setups of prior works, and implement strong LfS baselines by adopting shallow ConvNet
20 encoders and random shift augmentations. Surprisingly, we find that this modified LfS baseline can
21 achieve results competitive with prior works that leverage frozen pre-trained visual representations.
22 While our contributions are incremental in nature, we believe that our work contains must-know
23 insights for anyone working on pre-trained representations for motor control.

24 We evaluate our approach across **3** task domains (Adroit [21], DMControl [22], PixMC [3]) and **3**
25 algorithm classes: imitation learning (behavior cloning), on-policy RL (PPO [23]), and off-policy

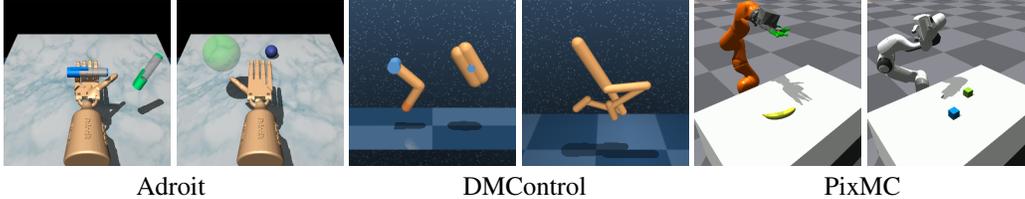


Figure 2: **Tasks.** We consider challenging and diverse visuo-motor control tasks spanning **3** domains: Adroit (dexterous manipulation), DMControl (locomotion, manipulation), and PixMC (manipulation).

26 RL (DrQ-v2 [2]). Our carefully designed LfS baseline is competitive with use of frozen pre-trained
 27 representations in all settings, and in some cases even outperforms them. We remain optimistic that
 28 pre-trained representations will play an important and increasingly larger role in visuo-motor control.
 29 At the same time, we believe that setting a simple yet strong baseline will help benchmark progress
 30 in this area. We conjecture that current benchmark tasks are not well suited to reap the benefits of
 31 pre-trained representations, since they do not require any visual generalization. As the community
 32 builds better benchmarks and harder tasks that require both visual and policy generalization, we
 33 conjecture that pre-trained representations will play an increasingly important role. *We are committed*
 34 *to releasing all of our proposed LfS baselines to the public.*

35 2 Experiments

36 Comparing two *paradigms* is difficult, and comparing LfS with pre-trained representations is no
 37 exception. To help narrow our scope, we focus on *representative methods* from each paradigm:
 38 a simple LfS method using data augmentation and a shallow ConvNet, and three **frozen** visual
 39 representations trained on large-scale out-of-domain vision datasets (PVR [1], MVP [3], R3M [13]).
 40 We choose to freeze the pre-trained representations to be consistent with prior work. All three
 41 pre-trained representations that we consider have been shown to outperform common representations
 42 such as supervised learning and MoCo-v2 [6] pretraining on ImageNet [16].

43 We propose a set of strong LfS baselines that span **3** classes of algorithms: imitation learning
 44 (behavior cloning), on-policy RL (PPO [23]), and off-policy RL (DrQ-v2 [2]), and consider a total
 45 of **15** tasks across **3** domains: Adroit [21] (dexterous manipulation; 2 tasks), DMControl [22]
 46 (locomotion and control; 5 tasks), and PixMC [3] (robotic manipulation; 8 tasks). Figure 2 shows
 47 sample tasks from each domain. We base our experiments on the public implementations of PVR,
 48 MVP, and DrQ-v2, and *meticulously follow their respective experimental setups*. We summarize our
 49 experiments setup as follows:

- 50 • **Behavior Cloning (BC).** We consider two domains – Adroit and DMControl – used in PVR.
 51 Observations are 256×256 RGB images (center-cropped to 224×224) with no access to pro-
 52 prioceptive information. Policies are trained with BC on 100 demonstrations per task; Adroit
 53 demonstrations are generated by oracle DAPG [21] policies as in PVR, and DMControl demon-
 54 strations are generated by oracle TD-MPC [24] policies. The original LfS baseline in PVR
 55 uses a shallow ConvNet encoder. Our improved LfS baseline additionally uses random shift
 56 augmentation [20, 2] during learning, and we refer to this baseline as *LfS (+aug)*. Data augmen-
 57 tation is relatively underexplored in BC literature, but we find that it works surprisingly well. In
 58 addition to PVR, we also compare with frozen MVP and R3M representations. Consistent with
 59 the experimental setup in PVR, we measure the policy performance with success rate in case of
 60 Adroit, and episode returns in case of DMControl. The policies are evaluated every two epochs,
 61 and we report the average performance over the three best epochs over the course of learning.
- 62 • **On-policy RL.** We reproduce the results of MVP on their proposed PixMC robotic manipulation
 63 benchmark. Observations are 224×224 RGB images and also include proprioceptive information.
 64 The original LfS baseline uses a small ViT [25] encoder. We propose *two* improved LfS baselines
 65 for this setting: (1) an LfS baseline that uses a shallow ConvNet encoder and *no* data augmentation,
 66 referred to as *LfS*, and (2) an LfS baseline that additionally applies random shift augmentation in
 67 critic learning, referred to as *LfS (+aug)*. Following prior work [4, 26], we do not augment value

68 targets. In addition to (frozen) MVP, we also compare with frozen PVR and R3M representations.
 69 Following the setup in MVP, we use success rate of the policy as the metric for comparison.

70 • **Off-policy RL.** We reproduce the results of state-of-the-art LfS method DrQ-v2 on the same
 71 DMControl tasks as used in PVR. Observations are 84×84 RGB images with no access to
 72 proprioceptive information; we upsample observations to 224×224 when using pre-trained
 73 representations. DrQ-v2 uses a shallow ConvNet encoder and random shift augmentation by
 74 default, and we refer to this baseline as *LfS (+aug)*. We compare DrQ-v2 to two alternatives:
 75 (1) not using data augmentation (simply denoted *LfS*), and (2) removing data augmentation **and**
 76 additionally replacing the LfS encoder with a frozen pre-trained representation, denoted by their
 77 representation names (PVR, R3M, MVP) respectively. Following prior work on DMControl, we
 78 use (normalized) return as the metric for comparison.

79 2.1 Results

80 We summarize our key findings as follows:

- 81 • **Performance comparison.** Our proposed Learning-from-Scratch (LfS) baselines are competitive
 82 with (and in some cases outperform) recent frozen pre-trained representations for visuo-motor
 83 control across a variety of algorithms and domains; see Figure 1 and Table 2. This indicates that,
 84 while pre-trained representations have the potential to replace the LfS paradigm in the future,
 85 under the set of most widely used metrics, they have yet to exceed the representational power of a
 86 *well designed* LfS method on standard benchmarks.
- 87 • **No free lunch – yet.** Our results indicate that the efficacy of a frozen pre-trained representation is
 88 both *task-dependent* (see Figure 3) and *algorithm-dependent* (see Figure 1): R3M outperforms
 89 other pre-trained representations on both Adroit and DMControl using BC, but performs poorly
 90 on DMControl using RL. Likewise, MVP performs well on PixelMC, but comparably worse in the
 91 two BC domains. Even within a visually consistent benchmark (PixMC), no single representation
 92 comes out on top. In contrast, LfS generally produces consistent results across all settings,
 93 presumably due to learning from task-specific data; this hypothesis is supported by prior work on
 94 in-domain finetuning of pretrained representations [27, 14, 28].

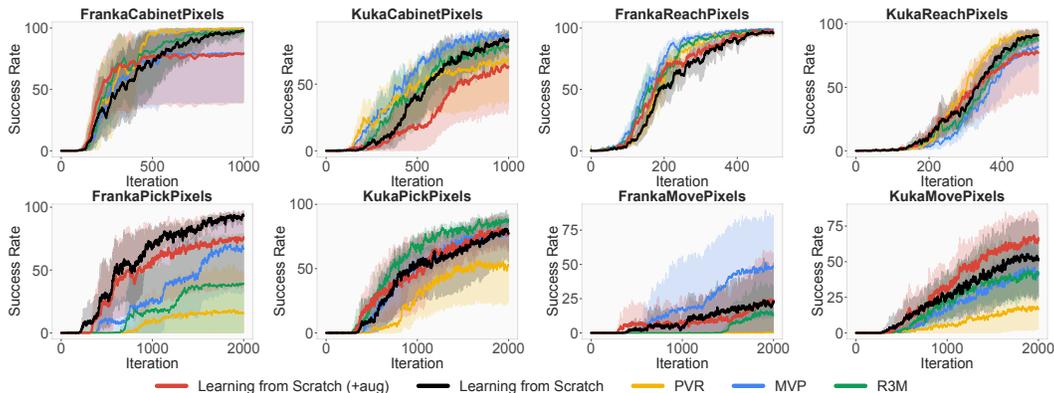


Figure 3: **PixMC benchmark.** PPO learning curves on the 8 robotic manipulation tasks from PixMC [3]. LfS performs comparably to pre-trained representations on most tasks. Averaged across 5 seeds.

Table 1: **Wall-time** of methods learning from scratch vs. using a pre-trained visual representation. For the latter, we report $\min\{\text{PVR}, \text{MVP}, \text{R3M}\}$ for a fair comparison. While LfS generally leads to better downstream task performance, using a **frozen** pre-trained representation can reduce computational cost substantially, especially during the training process. \downarrow Lower is better.

Method \ Setting	Behavior Cloning		Reinforcement Learning			
	Training (s/iteration)	Inference (s/episode)	s/1k frames	s/iteration		
LfS (+aug)	0.263	0.270	1.61	3.81	10.20	19.40
Fastest pre-training	0.003	0.006	2.66	11.00	13.00	11.90

95 • **Computational cost.** Our results so far has focused entirely on downstream task performance,
 96 like success rate or return. However, frozen pre-trained representations already demonstrate
 97 significant gains along an often-neglected axis: *wall-time*. Training and inference speeds are
 98 shown in Table 1. We find that BC policy updates are an order of magnitude faster using frozen
 99 pre-trained representations compared to LfS, as we can embed and cache features for the entire
 100 dataset in a few forward passes. However, inference speed generally favors LfS due to their
 101 smaller visual backbones, which is particularly important for real-robot applications. Since RL
 102 training interleaves learning and inference (data collection), wall-times are more balanced in this
 103 setting. We do not factor in the cost of learning a pretrained representation, since it is a one-time
 104 cost, and the representations can be reused across tasks.

105 3 Related Work

106 **Pre-training.** Representation
 107 learning via supervised/self-
 108 supervised/unsupervised pre-training
 109 on large-scale datasets has emerged
 110 as a powerful paradigm in areas
 111 such as computer vision [5, 6, 7, 8]
 112 and natural language processing
 113 [9, 10, 11, 12], where large datasets
 114 are available. While pre-trained
 115 representations can be finetuned to
 116 solve various downstream tasks, it
 117 may be prohibitively expensive to do
 118 so, and representations are therefore

Table 2: **Imitation Learning.** Success rate (Adroit) and unnormalized return (DMControl) of LfS and our **best** result obtained with a pre-trained representation, *i.e.*, for each task we report $\max\{\text{PVR, MVP, R3M}\}$.

Task\Method	LfS	LfS (+aug)	Best pre-training
Pen	14.0±5.6	85.7±4.3	87.0±5.7
Relocate	8.0±1.4	55.0±10.1	34.2±2.2
Finger Spin	6.2±3.0	445.3±12.2	611.4±18.2
Reacher Hard	35.6±18.1	846.2±67.2	602.0±107.8
Cheetah Run	8.8±5.6	171.0±18.8	202.2±12.2
Walker Stand	147.3±7.9	311.7±55.1	309.3±22.7
Walker Walk	38.4±2.4	111.4±27.5	80.6±5.3

119 commonly used as-is, *i.e.*, with *frozen* weights. We reflect on recent progress and challenges when
 120 leveraging pre-trained visual representations for control, which is an emerging and comparably
 121 underexplored application area of such representations.

122 **Pretrained representations for control.** Multiple works have explored learning control policies
 123 with visual representations pre-trained on large external datasets [29, 1, 13, 3, 27, 14, 15, 28]. In
 124 particular, PVR [1] and R3M [13] propose to learn policies by behavior cloning using pre-trained
 125 representations; PVR fuses features from several layers of a ResNet50 learned by MoCo [6], and R3M
 126 [13] learn a representation using a time-contrastive objective on ego-centric human videos. MVP
 127 [3] learn a policy with PPO and use a pre-trained visual encoder for feature extraction in addition to
 128 proprioceptive state information; the pre-trained representation is an MAE [30] trained on frames
 129 from diverse human videos. Despite encouraging results in leveraging pre-trained representations for
 130 control, we show that LfS remains competitive with (frozen) pre-trained representations at this time.

131 4 Discussion

132 We have shown that a carefully designed LfS baseline is competitive with frozen pre-trained repre-
 133 sentations across a variety of algorithm classes and domains. While this is the current conclusion, we
 134 remain optimistic that results will be skewed in favor of pre-trained representations as the paradigm
 135 matures. At present, we find that the main benefit of a frozen pretrained representation is the reduced
 136 training cost that comes with its *universality* – a single representation can be reused across tasks.
 137 Bridging the performance gap while maintaining universality will thus be critical to the adoption
 138 of this new paradigm. Recent works show that pretrained representations benefit from finetuning
 139 on task-specific data [27, 14, 28], combining elements of pre-training and LfS. However, finetuning
 140 large visual backbones present optimization challenges (*e.g.*, instability and catastrophic forgetting),
 141 and can be costly. We encourage further research in these directions, and hope that our strong LfS
 142 baselines will help accurately benchmark progress in this area.

References

- 143
- 144 [1] S. Parisi, A. Rajeswaran, S. Purushwalkam, and A. K. Gupta. The unsurprising effectiveness of
145 pre-trained vision models for control. In *ICML*, 2022.
- 146 [2] D. Yarats, I. Kostrikov, and R. Fergus. Image augmentation is all you need: Regularizing deep
147 reinforcement learning from pixels. In *International Conference on Learning Representations*,
148 2021. URL <https://openreview.net/forum?id=GY6-6sTvGaf>.
- 149 [3] T. Xiao, I. Radosavovic, T. Darrell, and J. Malik. Masked visual pre-training for motor control.
150 *ArXiv*, abs/2203.06173, 2022.
- 151 [4] N. Hansen, H. Su, and X. Wang. Stabilizing deep q-learning with convnets and vision trans-
152 formers under data augmentation. In *NeurIPS*, 2021.
- 153 [5] C. Doersch, A. K. Gupta, and A. A. Efros. Unsupervised visual representation learning by
154 context prediction. *2015 IEEE International Conference on Computer Vision (ICCV)*, pages
155 1422–1430, 2015.
- 156 [6] K. He, H. Fan, Y. Wu, S. Xie, and R. B. Girshick. Momentum contrast for unsupervised
157 visual representation learning. *2020 IEEE/CVF Conference on Computer Vision and Pattern
158 Recognition (CVPR)*, pages 9726–9735, 2020.
- 159 [7] A. van den Oord, Y. Li, and O. Vinyals. Representation learning with contrastive predictive
160 coding. *ArXiv*, abs/1807.03748, 2018.
- 161 [8] J.-B. Alayrac, J. Donahue, P. Luc, A. Miech, I. Barr, Y. Hasson, K. Lenc, A. Mensch, K. Millican,
162 M. Reynolds, et al. Flamingo: a visual language model for few-shot learning. *arXiv preprint
163 arXiv:2204.14198*, 2022.
- 164 [9] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional
165 transformers for language understanding. *ArXiv*, abs/1810.04805, 2019.
- 166 [10] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam,
167 G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. J. Henighan, R. Child,
168 A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray,
169 B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei.
170 Language models are few-shot learners. *ArXiv*, abs/2005.14165, 2020.
- 171 [11] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell,
172 P. Mishkin, J. Clark, et al. Learning transferable visual models from natural language supervision.
173 In *International Conference on Machine Learning*, pages 8748–8763, 2021.
- 174 [12] A. Chowdhery, S. Narang, J. Devlin, M. Bosma, G. Mishra, A. Roberts, P. Barham, H. W.
175 Chung, C. Sutton, S. Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv
176 preprint arXiv:2204.02311*, 2022.
- 177 [13] S. Nair, A. Rajeswaran, V. Kumar, C. Finn, and A. Gupta. R3m: A universal visual representation
178 for robot manipulation. *ArXiv*, abs/2203.12601, 2022.
- 179 [14] Y. Ze, N. Hansen, Y. Chen, M. Jain, and X. Wang. Visual reinforcement learning with self-
180 supervised 3d representations. *arXiv preprint arXiv:2210.07241*, 2022.
- 181 [15] Z. Yuan, Z. Xue, B. Yuan, X. Wang, Y. Wu, Y. Gao, and H. Xu. Pre-trained image encoder for
182 generalizable visual reinforcement learning. In *First Workshop on Pre-training: Perspectives,
183 Pitfalls, and Paths Forward at ICML 2022*, 2022. URL [https://openreview.net/forum?
184 id=E-0zNz5J5BM](https://openreview.net/forum?id=E-0zNz5J5BM).
- 185 [16] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy,
186 A. Khosla, M. S. Bernstein, A. C. Berg, and L. Fei-Fei. Imagenet large scale visual recognition
187 challenge. *International Journal of Computer Vision*, 115:211–252, 2015.

- 188 [17] K. Grauman, A. Westbury, E. Byrne, Z. Chavis, A. Furnari, R. Girdhar, J. Hamburger, H. Jiang,
189 M. Liu, X. Liu, et al. Ego4d: Around the world in 3,000 hours of egocentric video. In
190 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
191 18995–19012, 2022.
- 192 [18] A. Srinivas, M. Laskin, and P. Abbeel. Curl: Contrastive unsupervised representations for
193 reinforcement learning. In *ICML*, 2020.
- 194 [19] M. Schwarzer, A. Anand, R. Goel, R. D. Hjelm, A. C. Courville, and P. Bachman. Data-efficient
195 reinforcement learning with self-predictive representations. In *ICLR*, 2021.
- 196 [20] I. Kostrikov, D. Yarats, and R. Fergus. Image augmentation is all you need: Regularizing deep
197 reinforcement learning from pixels. *ArXiv*, abs/2004.13649, 2021.
- 198 [21] A. Rajeswaran, V. Kumar, A. Gupta, G. Vezzani, J. Schulman, E. Todorov, and S. Levine. Learn-
199 ing Complex Dexterous Manipulation with Deep Reinforcement Learning and Demonstrations.
200 In *Proceedings of Robotics: Science and Systems (RSS)*, 2018.
- 201 [22] Y. Tassa, Y. Doron, A. Muldal, T. Erez, Y. Li, D. de Las Casas, D. Budden, A. Abdolmaleki,
202 et al. Deepmind control suite. Technical report, DeepMind, 2018.
- 203 [23] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization
204 algorithms. *ArXiv*, abs/1707.06347, 2017.
- 205 [24] N. Hansen, X. Wang, and H. Su. Temporal difference learning for model predictive control. In
206 *ICML*, 2022.
- 207 [25] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani,
208 M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby. An image is worth 16x16
209 words: Transformers for image recognition at scale. *ArXiv*, abs/2010.11929, 2021.
- 210 [26] R. Raileanu, M. Goldstein, D. Yarats, I. Kostrikov, and R. Fergus. Automatic data augmentation
211 for generalization in deep reinforcement learning. *arXiv preprint arXiv:2006.12862*, 2020.
- 212 [27] C. Wang, X. Luo, K. W. Ross, and D. Li. Vr13: A data-driven framework for visual deep
213 reinforcement learning. *ArXiv*, abs/2202.10324, 2022.
- 214 [28] Y. Xu, N. Hansen, Z. Wang, Y.-C. Chan, H. Su, and Z. Tu. On the feasibility of cross-task
215 transfer with model-based reinforcement learning. *arXiv preprint arXiv:2210.10763*, 2022.
- 216 [29] R. Shah and V. Kumar. Rrl: Resnet as representation for reinforcement learning. *ArXiv*,
217 abs/2107.03380, 2021.
- 218 [30] K. He, X. Chen, S. Xie, Y. Li, P. Doll’ar, and R. B. Girshick. Masked autoencoders are scalable
219 vision learners. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition*
220 *(CVPR)*, pages 15979–15988, 2022.