
Masked Visual Pre-training for Motor Control

Tete Xiao^{*1} Ilija Radosavovic^{*1} Trevor Darrell^{†1} Jitendra Malik^{†1}

Abstract

This paper shows that self-supervised visual pre-training from real-world images is effective for learning motor control tasks from pixels. We first train the visual representations by masked modeling of natural images. We then freeze the visual encoder and train neural network controllers on top with reinforcement learning. We do not perform any task-specific fine-tuning of the encoder; the same visual representations are used for all motor control tasks. To the best of our knowledge, this is the first self-supervised model to exploit real-world images at scale for motor control. To accelerate progress in learning from pixels, we contribute a benchmark suite of hand-designed tasks varying in movements, scenes, and robots. Without relying on labels, state-estimation, or expert demonstrations, we consistently outperform supervised encoders by up to 80% absolute success rate, sometimes even matching the oracle state performance. We also find that in-the-wild images, e.g., from YouTube or Egocentric videos, lead to better visual representations for various manipulation tasks than ImageNet images.

1. Introduction

The last decade of machine learning has been powered by learning representations with large neural networks and augmenting them with a relatively small amount of domain knowledge where appropriate. This paradigm has led to substantial progress across a range of domains. Examples include visual recognition (Girshick et al., 2014; He et al., 2017), natural language (Radford et al., 2018; Devlin et al., 2019; Radford et al., 2019; Brown et al., 2020), and audio (Van Den Oord et al., 2016). And the trend continues. Motor control, however, remains a notable exception.

^{*,†}Equal contribution ¹University of California, Berkeley. Correspondence to: Tete Xiao <txiao@eecs.berkeley.edu>, Ilija Radosavovic <ilija@berkeley.edu>.

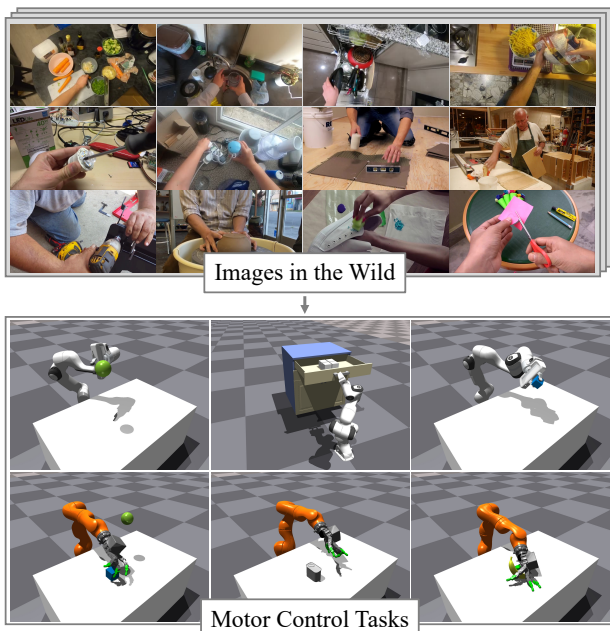


Figure 1. We explore learning complex motor control from pixels. We show that we are able to solve a range of motor control tasks with variations in robots, scenes, and objects. This is enabled by learning rich visual representations from real-world images with masked modeling. Videos available on the [project page](#).

In this paper, we show that self-supervised visual pre-training on real-world images is effective for learning motor control tasks from pixels. These self-supervised representations consistently outperform supervised representations.

Consider tasks shown in Figure 1 (bottom). The required movement types vary from simple reaching to object interactions. We also see variations in robots, scene configurations, and objects. Control inputs are high-dimensional and difficult to search (e.g., 23 DoF robot with a multi-finger hand). We explore learning complex tasks such as these from high-dimensional pixel observations.

To tackle this setting, we use the neural network architecture shown in Figure 2b. Our network encodes the input image using a high-capacity visual encoder (Dosovitskiy et al., 2020) and combines it with proprioceptive information to obtain an embedding. A light-weight neural network controller takes in the embedding and predicts actions. The whole system can be trained end-to-end.

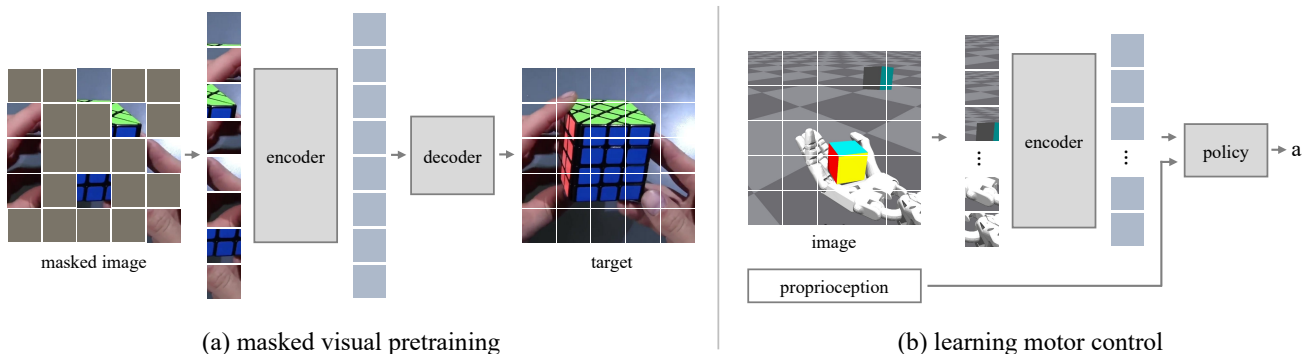


Figure 2. **Masked visual pre-training for motor control.** *Left:* We first *pre-train* visual representations using *self-supervision* through masked image modeling (He et al., 2021) from *real-world* images. *Right:* We then *freeze* the image encoder and train task-specific controllers on top with reinforcement learning (RL). The *same* visual representations are used for all motor control tasks.

Indeed, this design is akin to architectures typically used in approaches that learn control policies end-to-end with RL, e.g., Levine et al. (2016). While conceptually appealing, the latter has two main challenges in practice. First, training is computationally expensive and has poor sample complexity (especially with high-dimensional inputs and actions). Second, the learned solutions typically overfit to the setting at hand and thus do not generalize to new scenes and objects.

One way to offset the high sample complexity of end-to-end RL is to employ auxiliary objectives (Jaderberg et al., 2017; Oord et al., 2018; Yarats et al., 2019; Srinivas et al., 2020). For example, Srinivas et al. (2020) show excellent performance in vision-based RL by using contrastive learning with data augmentations. However, such representations are still trained using only environment-specific experience.

The key aspect of our approach is in how we train the visual representations. We do *not* train the visual encoder while learning specific motor control tasks. Instead, we *pre-train* the visual encoder by *self-supervision* from *natural images* (Figure 1 & 2). We learn the visual representation by performing masked image modeling through the masked autoencoder (MAE) (He et al., 2021). Thanks to the Internet and ubiquitous portable cameras, we now have access to large amounts of unlabeled visual data for self-supervision. MAE does not require human labels or make strong assumptions about data distributions, e.g., centered objects or pre-defined augmentation invariances (Xiao et al., 2021b), making it an excellent framework for learning general visual representations from large collections of in-the-wild images.

Given the visual encoder, we train controllers on top with reinforcement learning (Schulman et al., 2017). We keep the visual representations frozen and do not perform *any* task-specific fine-tuning of the encoder; all motor control tasks use the *same* visual representations. We call our approach MVP (for **M**asked **V**isual **P**re-training for **M**otor **C**ontrol).

To accelerate future progress, we contribute PixMC, a new benchmark suite of hand-designed tasks with various movement types, scene configurations, and robots. We leverage a GPU-based simulator for fast simulation (Makoviychuk et al., 2021), provide reward functions, baselines, and multi-GPU implementation of learning algorithms from pixels.

We compare our self-supervised approach to baselines that follow the same architecture (Figure 2b) but use different visual representations. As an upper bound, we consider oracle hand-engineered states for solving a task (e.g., 3D poses and direction-to-goal vectors). We also compare our method to visual encoders trained by supervised learning on ImageNet (Deng et al., 2009), the choice of encoder in most vision tasks. We summarize our main results as follows:

- 1) We show that a *single* visual encoder pre-trained on real-world images can solve various motor control tasks *without* fine-tuning per-task, state estimation, or demonstrations.
- 2) Our *self-supervised* approach consistently outperforms *supervised* representations (up to 80% absolute success rate), and even matches the oracle performance in some cases.
- 3) We find that pre-training on *images in the wild*, e.g., from YouTube (Shan et al., 2020) or Egocentric (Damen et al., 2018; 2021) videos, works better for manipulation tasks than ImageNet (Deng et al., 2009) images.
- 4) We show that our visual representations *generalize* in various ways. For example, our visual encoder disentangles shape and color and is able to handle a range of different object geometries and configurations.

We encourage researchers to evaluate visual representations not only on downstream vision tasks but also on motor control tasks. We believe that our work is a promising step in this direction and release the benchmark suite, pre-trained models, and the training code on the [project page](#).

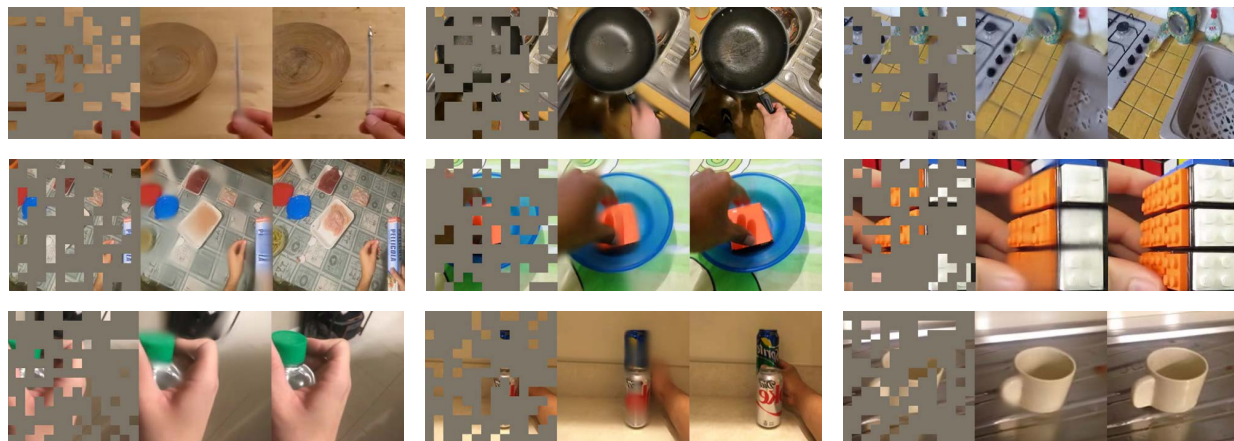


Figure 3. **Example reconstructions.** For each triplet from left to right: the masked image, the reconstructed image, the ground-truth target. We observe that the autoencoder learns color, shape, and object affordance using self-supervision from in-the-wild images.

2. Masked Visual Pre-training for MC

2.1. Masked Visual Pre-training

Thanks to the Internet and portable camera devices (e.g., phones, glasses, etc.), we now have access to large amounts of video data to learn from in various ways. We leverage such data for learning *visual representations*. Specifically, we use images from the egocentric Epic Kitchens dataset (Damen et al., 2018; 2021), the YouTube 100 Days of Hands dataset (Shan et al., 2020), and the crowd-sourced Something-Something dataset (Goyal et al., 2017b). Combined, these sources yield a collection of $\sim 700\text{K}$ images that we refer to as the Human-Object Interaction dataset (HOI). Note that we do *not* exploit any human labels or temporal information even if it is possible. We also apply our approach using the ImageNet dataset (Deng et al., 2009) for controlled comparisons with supervised baselines.

With the data in hand, we must now formulate an appropriate self-supervised task. We adopt masked modeling as our self-supervision objective—specifically, we use masked autoencoder (MAE) (He et al., 2021). MAEs mask-out random patches of the input image and reconstruct the missing pixels with a Vision Transformer (ViT) (Dosovitskiy et al., 2020). During training, only unmasked patches are fed into the MAE; this strategy makes training more efficient. It is critical to train MAE with a high masking ratio (e.g., masking 75% of all patches) and use a heavy encoder with a light decoder. The most appealing property of MAE is its simplicity and minimal reliance on dataset-specific augmentation engineering; for example, it works well even with minimal data augmentations (center crop and color).

In Figure 3 we show example reconstructions for HOI images. We observe that the model learns about color, shape, and objects. Notice that the images are representative of everyday interactions making them well suited for our needs.

2.2. Learning Motor Control from Pixels

Given the pre-trained visual encoder, we now turn to learning motor control from pixels. We freeze the visual representations and use them for all downstream motor control tasks; we do not perform any task-specific fine-tuning of the image encoder. This design has two main benefits. First, it prevents the encoder from overfitting to the setting at hand and thus preserves general visual representations for learning new tasks. Second, it leads to considerable memory and run time savings since there is no need to back-propagate through the encoder. Freezing visual encoder makes using large vision models in the RL loop fast and feasible.

In Figure 2b, we show our architecture for learning motor control from pixels. We first extract a fixed-sized vector of image features using our pre-trained visual encoder. Notice that all of the image patches are passed through the encoder, unlike in the masked pre-training stage. We additionally compile proprioceptive robot information in the form of joint positions and velocities into a second vector. This proprioceptive information is readily available on real robot hardware. We concatenate these two vectors to obtain the input embedding for the neural network controller.

We then train task-specific motor control policies on top of this embedding with model-free reinforcement learning. Specifically, we use the proximal policy optimization (PPO) algorithm (Schulman et al., 2017). PPO is a state-of-the-art policy gradient method that has shown excellent performance on complex motor control tasks and successful transfer to real hardware (OpenAI et al., 2020; 2019). Our policy is a small multi-layer perceptron (MLP) network. In addition, we train a critic that has the same architecture as the policy using the same representations. The policy and the critic do not share weights.

	RLBENCH	ROBOSUITE	METAWORLD	OURS
SIMULATOR	COPPELIA	MUJoCo	MUJoCo	ISAACGYM
FAST				✓
#ARMS	1	8	1	2
#HANDS				✓
#TASKS	100	9	50	8
REWARDS		✓	✓	✓

Table 1. **Existing benchmarks.** Compared to existing benchmarks, ours features a unique combination of hand-designed tasks, dense rewards, and complex robots (e.g., multi-finger hands). Crucially, it leverages a fast simulator and provides distributed training for scaling learning-based motor control from pixel observations.

3. PixMC Benchmark

We construct a new benchmark suite of tasks for studying motor control from pixels, described in this section.

3.1. Motivation

While there exist a number of excellent benchmarks for motor control, e.g., DMC (Tassa et al., 2018), RLBench (James et al., 2020), Robosuite (Zhu et al., 2020), MetaWorld (Yu et al., 2020), they all fall short on one or more of our requirements. In particular, there is no benchmark suite for learning motor control algorithms that has high-resolution images, realistic robots, fast physics simulation, efficient training, and appropriate reward functions and metrics. To this end, we introduce a new benchmark suite for *Pixel Motor Control*, which we call PixMC. We compare the key aspects of PixMC to several existing benchmarks in Table 1.

3.2. Simulator

We leverage the recent NVIDIA IsaacGym simulator (Makoviychuk et al., 2021) to build our benchmark. The core design idea of IsaacGym is to perform simulation on a GPU in a shared context. IsaacGym allows for very fast training times. For example, we are able to train our oracle state-based models in ~ 12 minutes and our image models in ~ 5 hours (~ 8 million environment steps) on a single NVIDIA 2080 Ti GPU. This training speed is considerably faster than it would be in other simulators commonly used for motor control such as MuJoCo (Todorov et al., 2012). Rudin et al. (2021) have shown that sim-to-real transfer is feasible based on policies trained in IsaacGym.

3.3. Robots

PixMC includes two robot arms, a parallel jaw gripper, and a multi-finger hand combined as follows: **(1) Franka:** A Franka Emika robot commonly used for research. It has a 7-DoF arm with a 2-DoF gripper. **(2) Kuka with Allegro:** A Kuka LBR iiwa arm with 7 DoFs and a 4-finger Allegro hand with 16 DoFs (4 DoFs per finger), for 23 DoFs in total. For brevity, we refer to it as “Kuka.”

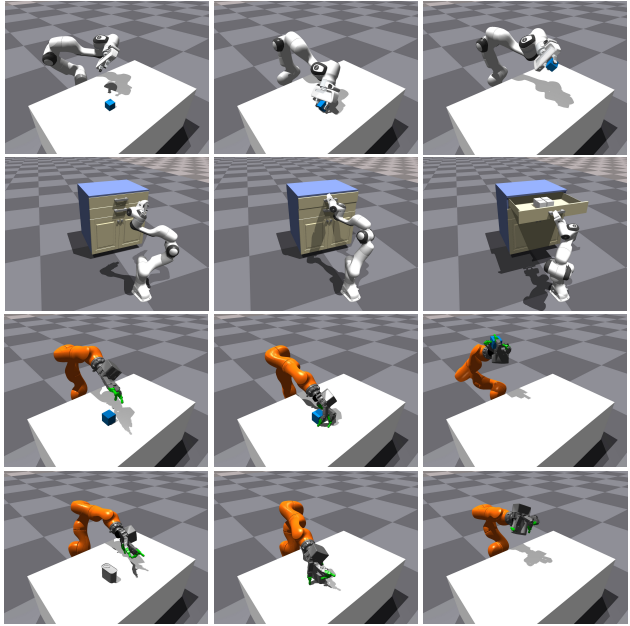


Figure 4. **Example tasks.** We show example trajectories for the Franka and Kuka tasks. See the [project page](#) for video examples.

3.4. Observations and Actions

Our benchmark supports rendering high-resolution pixel observations for each robot. For both the Franka and Kuka setups and each of the tasks, we use wrist-mounted cameras by default. The benchmark provides proprioceptive information for the robots, as well as hand-engineered states typically including 3D poses or relevant objects, goals, and their relations. All of our default settings use position control in joint angle space with a control frequency of 60Hz.

3.5. Tasks, Rewards, and Metrics

PixMC tasks include several movement types from basic reaching to interacting with objects. The objects in the environment vary in positions, scale, color, and shape. Figure 1 and 4 show a few example scenes. We hand-design task-specific dense reward functions for training RL policies. We define reward-independent success metrics that typically quantify the distance from the agent or an object to a specified goal location over sufficient time steps.

3.6. Distributed Training

The scarcity of GPU memory is a bottleneck for learning motor control from pixels. For our typical setup with 224×224 images, we can fit at most 256 environments on a single 2080 Ti GPU. We implement PPO with distributed training to support large batch sizes. Similar to data parallel training, we create a model replica per-GPU, collect rollouts on each GPU, and synchronize gradients across GPUs.

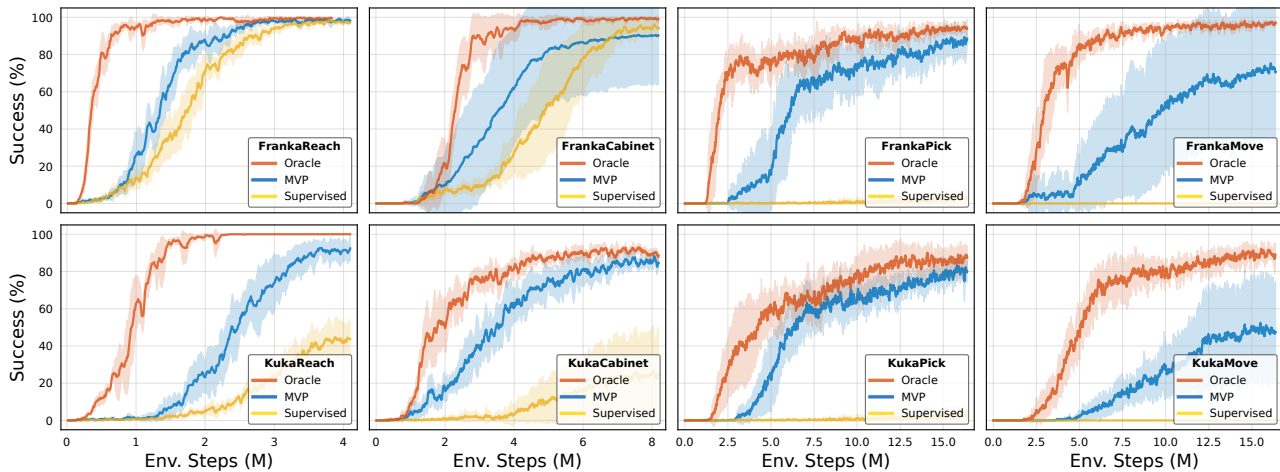


Figure 5. **Sample complexity.** We plot the success rate as a function of environment steps on the 8 PixMC tasks. Each task uses either the Franka arm with a parallel gripper or the Kuka arm with a multi-finger hand. The MVP approach significantly outperforms the supervised baseline on 7 tasks and closely matches the oracle state model (considered the upper bound of RL) on 5 tasks at convergence. The result shows that self-supervised pre-training markedly improves representation quality for motor control tasks.

4. Experimental Setup

Data for pre-training. We consider two kinds of pre-training data: ImageNet (Deng et al., 2009) and a joint Human-Object Interaction (HOI) dataset. We construct the HOI data by combining Epic-Kitchens (Damen et al., 2018; 2021), Something-Something (Goyal et al., 2017b), and 100 Days of Hands (100-DOH) (Shan et al., 2020). To build HOI, we sample frames from Epic-Kitchens and Something-Something at 1fps and 0.3fps, respectively. This yields 700k images including 100k from 100-DOH.

Encoder. The image encoder follows standard ViT architecture (Dosovitskiy et al., 2020). ViT partitions an image into patches and linearly projects each patch into features, followed by standard Transformer blocks (Vaswani et al., 2017; Wang et al., 2019). We use the ViT-Small model with a 16×16 patch size, 384 hidden size, 6 attention heads, an MLP multiplier of 4, and 12 blocks. The model runs at 4.6 gigaflops for input images of 224×224 , approximately $1.2 \times$ as many as the ResNet-50 (He et al., 2016) model.

We pre-train supervised and self-supervised variants of the ViT model. For the supervised model, we use the recipe in (Xiao et al., 2021a) and train on the ImageNet dataset for 400 epochs. We use the MAE framework for the self-supervised counterpart (He et al., 2021). We use an auxiliary dummy classification token in the MAE for downstream finetuning and transfer (He et al., 2021). We use a crop ratio of $[0.2, 1.0]$ for ImageNet and $[0.1, 0.75]$ for HOI, due to the larger width-over-height aspect ratio of HOI images. We train the MAE models for 1600 epochs on 16 GPUs for both HOI and ImageNet datasets.

Controller. The controller is a simple MLP with each hidden layer followed by a SeLU (Klambauer et al., 2017) activation function. We use a four-layer MLP with hidden layers of size $[256, 128, 64]$ for all tasks, following (Makoviychuk et al., 2021). The (dummy) classification token of the ViT encoder yields the image features and a linear layer projects the features to 128 dimensions. The controller takes in the linearly-projected (128-d) proprioceptive state of the robot along with the projected image features. The controller outputs delta joint angles.

Training with RL. We freeze the visual encoder throughout the entire training horizon. We train for 500 iterations for reach, 1000 iterations for cabinet, and 2000 iterations for pick and relocate, respectively. In each iteration, we collect samples from 256 environments which have 32 steps each. We train for 10 epochs on these collected samples per iteration. We compose 4 minibatches per epoch, i.e., 4 gradient updates, leading to a minibatch size of 2048 per gradient update. We choose this configuration because it maximizes the memory on a single NVIDIA 2080 Ti GPU. In all experiments we train with a cosine learning rate decay schedule (Goyal et al., 2017a). To reduce randomness in the RL experiments (Agarwal et al., 2021), for each task and model we search for the best learning rate in $\{0.0005, 0.001, 0.0015\}$ with 5 seeds per learning rate (15 runs for each task and model). We always report the performance yielded by the best learning rate aggregated over seeds unless otherwise specified. Other hyperparams use defaults: Adam optimizer with $\beta_1 = 0.9$ and $\beta_2 = 0.999$, gradient norm of 1, initial noise standard deviation of 1.0.

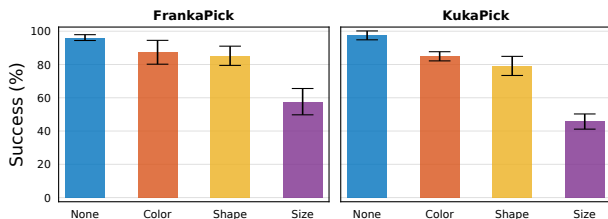
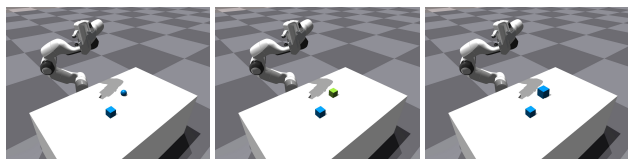


Figure 6. Robustness to distractors. The robots are trained to pick up a blue box of 4.5cm side length. At *test time*, we add a distractor object differing from the training object in terms of color (blue vs. green), shape (cube vs. sphere), or size (4.5cm vs. 6cm), shown at the top. MVP maintains high success rates for color and shape. The model is less sensitive to size likely due to the scale ambiguity from single first-person camera setup.

5. Experimental Results

5.1. Sample Complexity

Figure 5 shows success rates over training on 8 challenging tasks from PixMC. We consider the oracle state model (i.e., position, orientation, and velocity of the object, goal and robot in world-coordinate system, which is difficult to estimate in real-world settings) as the upper bound of RL. MVP significantly outperforms the supervised baseline on 7 out of 8 tasks and matches the baseline on the 8th task. At convergence, MVP closely matches the oracle state model on 5 tasks. The supervised baseline is flat at zero success rate on the pick and move tasks with both robots; MVP rivals the oracle on the pick task and achieves high success rate on the relocate task. These results show that self-supervised pre-training markedly improves representation quality for motor control tasks.

5.2. Generalization

We design various experiments on the pick task with both the Franka and Kuka arms to demonstrate the degree to which MVP is able to generalize.

Robustness to distractors. In the pick task the robots are trained to pick up a blue box of 4.5cm side length. In this experiment we add a distractor object that differs from the training object in terms of color, shape, or scale, *at testing time*. The models are not retrained for new testing configurations. A robust model from pixels should pick up the object used for training. Specifically, we have 1) a *green* box of 4.5cm side length (color distractor); 2) a blue *sphere* of 4.5cm diameter (shape distractor); and 3) a blue

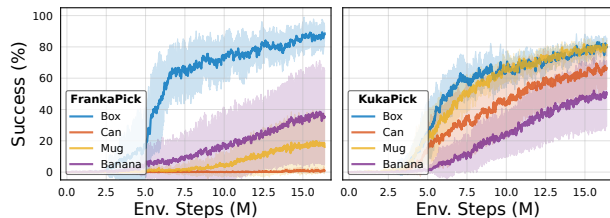


Figure 7. Generalization to objects of various geometries. We import three additional objects (i.e., can, mug, and banana) from the YCB dataset and re-train the controller to pick up the individual object category. The Kuka robot with the Allegro hand can pick up all objects with $\sim 50\%$ success rate.

box of 6cm side length (scale distractor). Figure 6 shows the results of our MVP model. Color and shape distractors only marginally decrease the success rate, implying that MVP is able to recognize the color and shape of objects. The model, however, is less sensitive to scale variation as the 50% success rate suggests that the distractor or the original box is picked up by chance. We believe it is due to scale ambiguity from single first-person camera setup. Note that the oracle state model would not be able to function in these testing cases due to changes in state dimensions.

Generality of the framework. We import various objects from the YCB dataset (Calli et al., 2015)—box, can, mug, and banana—for the pick task and re-train the model for each individual object. Figure 7 visualizes the experiment setup and shows the results. The Kuka robot with the Allegro hand can pick up all of the objects with at least a 50% success rate. This shows MVP as a framework can generalize to objects of different geometries and the multi-finger hand’s strength in object manipulation.

5.3. Ablations

Pre-training data. We train MVP on HOI and ImageNet data, respectively. Figure 8 shows the results. MVP trained on HOI data outperforms the counterpart trained on ImageNet data on 7 out of 8 tasks. Whereas ImageNet is dominated by images of animals and objects, HOI contains many images demonstrating object manipulation from a first-person camera view. We hypothesize that this difference is why HOI is empirically the superior choice for motor control tasks.

Random features. We compare our MVP with a randomly initialized and frozen image encoder. The random model

Masked Visual Pre-training for Motor Control

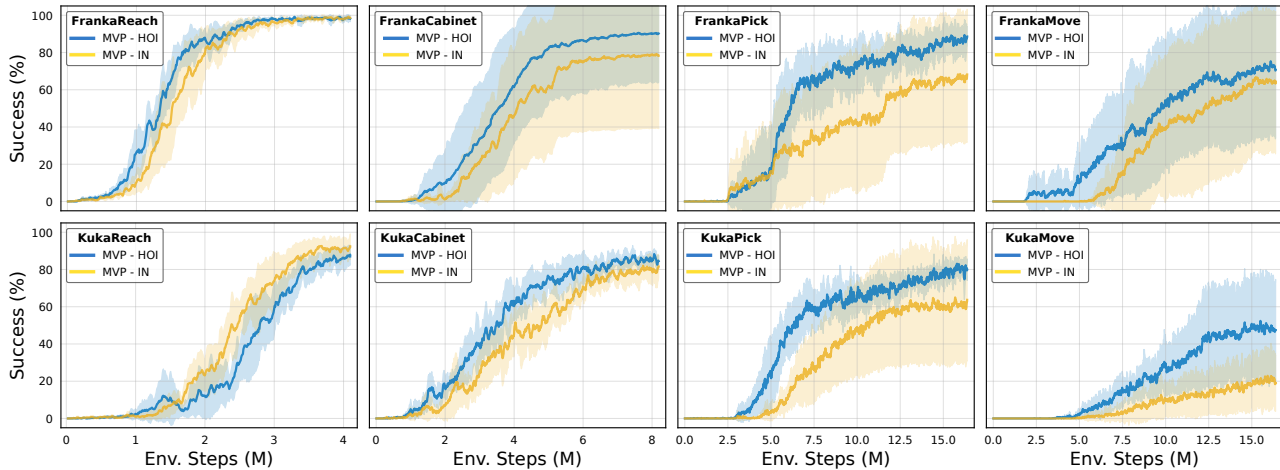


Figure 8. **Pre-training data: HOI vs. ImageNet.** We compare the performance of our approach when using HOI and ImageNet images as the source of pre-training data. Overall, we find that the representations learned from HOI data perform better on motor control tasks.

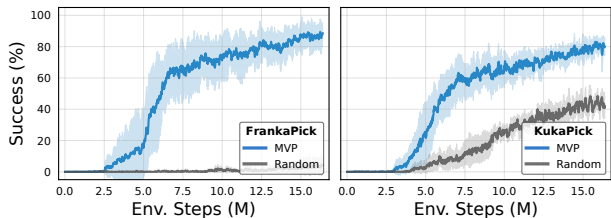


Figure 9. **Random features.** We compare our learnt representations to a random features baseline. We use the same visual encoder, initialize it randomly, and freeze. The random model fails on 6 out of 8 PixMC tasks (0 success rate). Here we show one task that fails (FrankaPick) and one that succeeds (KukaPick).

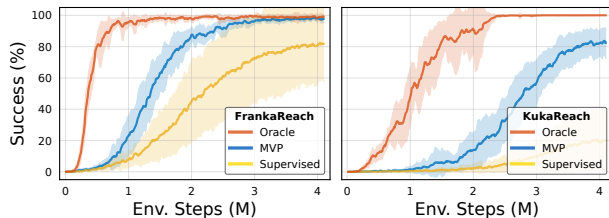


Figure 10. **Learning rate and seed stability.** For each model, we train 15 instances of the model with 3 learning rates and 5 seeds. We show the results on tasks where the supervised baseline achieves nontrivial performance. MVP significantly improves stability over the supervised model.

yields flat-zero or close-to-zero success rates on 6 out of 8 tasks. This indicates that random features are insufficient to solve complex motor control tasks. We show the results of one successful and one failed task in Figure 9.

Stability. We characterize how sensitive different models are to changes in learning rate and random seed. Representations of high quality should be less sensitive and yield consistently-high performance. Figure 10 shows the reach task results over 3 learning rates (0.0005, 0.001, 0.0015) and 5 seeds for a total of 15 runs for each model. Although the supervised baseline with the best choice of learning rate performs close to MVP in Franka reach and has over 40% success rate in Kuka reach (see Figure 5), it exhibits much worse stability across a range of learning rates. Our MVP, in contrast, shows good stability in the learning rate hyperparameter, a sign of superior optimizability in training.

Unfreezing encoders. We experiment with training the encoder end-to-end with RL in downstream motor control tasks. Unfreezing the encoder significantly increases GPU

compute and memory consumption. To maintain the number of environments, we increase the number of GPUs used from 1 to 8 via our distributed training pipeline. We choose the Franka pick task as the testbed and run 3 runs with different seeds per model. We test with two initialization: (1) initializing the visual encoder randomly, or (2) initializing with the pre-trained MAE weights. We observed unstable training (the loss goes to NaN), and we decreased the learning rate until training successfully completed. Still, both models yielded flat zero success rate on the task at all seeds. This result is somewhat counter-intuitive as one would expect end-to-end training should yield better results than training on frozen representations. We conjecture that it may be due to (a) the RL signal being unstable and hard to tune; (b) noisy gradients from the RL objective interfering with pre-trained visual representation; or (c) a high capacity vision model like ViT requiring significantly more samples and environments to train end-to-end with RL. Freezing the visual encoder as in MVP preserves the quality of visual representations while yielding faster RL training.

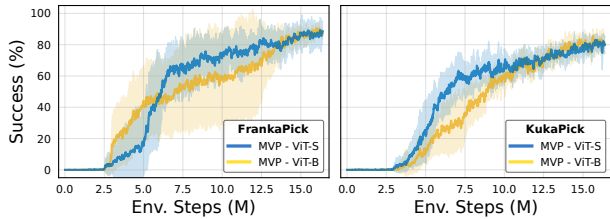


Figure 11. Larger encoders. We pre-train a ViT-Base model (18Gflops) and use the representations for the pick task. We do not observe clear gains from preliminary model scaling and believe that scaling data and model size is an exciting area for future work.

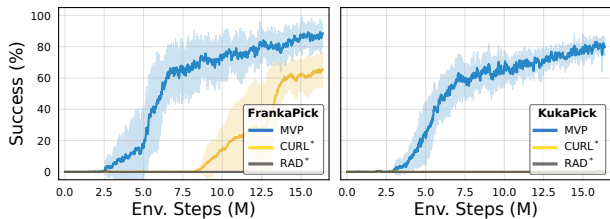


Figure 12. CURL and RAD comparisons. We compare MVP to environment training with our reimplementation of CURL and RAD (denoted with star; see text for details). We observe that, for a fixed number of steps, MVP outperforms both while being less computationally expensive to train (1 vs. 8 GPU training).

Larger encoders. We pre-train a ViT-Base encoder (18 gigaflops) and conduct a preliminary transfer study on the Franka/Kuka pick tasks. Figure 11 shows the results. We observe that the larger encoder does not improve performance. A larger encoder potentially requires more data and/or a different training recipe. Overall, scaling data and models in the context of self-supervised representations for motor control remains an exciting area for future work.

5.4. Additional Comparisons

CURL and RAD comparisons. We compare MVP to two state-of-the-art methods that train the vision encoder with environment data: CURL (Srinivas et al., 2020) and RAD (Laskin et al., 2020). Due to differences in their original settings and ours, e.g., small ConvNet vs. large ViT, we compare to our reimplementations (denoted with star). In particular, we adopt PPO, ViT-Small visual encoder, and MoCo-v3 (Chen et al., 2021b) data augmentation recipe. We observe that, for a fixed number of steps, MVP outperforms both baselines while being less computationally expensive (1 vs. 8 GPUs). We note that environment training with enough data might yield better performance than our image pre-training. Indeed, we do not see our approach as the ultimate answer for any one setting but rather as a solid baseline, or a starting point, for many varying settings.

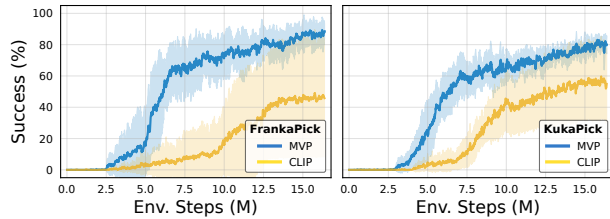


Figure 13. CLIP comparisons. We compare our visual representations to the CLIP visual encoder trained with large scale language supervision. The results show a promising signal that self-supervised representations can outperform CLIP encoders.

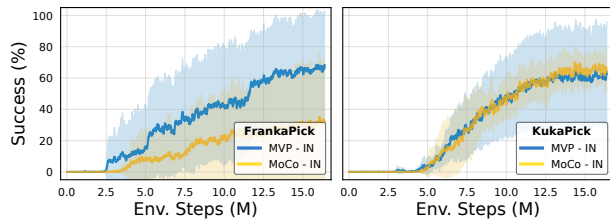


Figure 14. MoCo-v3 comparisons. We compare MAE, used in MVP, to the MoCo-v3 alternative for self-supervised pre-training on *ImageNet*. We see that MoCo-v3 can achieve non-trivial performance, showing the generality of our approach. However, in contrast to MAE, it may be harder to adapt to in-the-wild images.

CLIP comparisons. Next, we experiment with substituting a pre-trained CLIP visual encoder (Radford et al., 2021) in place of our MAE encoders. The CLIP encoder is trained using 400M labeled text-image pairs and has shown excellent performance across a wide range of visual tasks. We opt to use the ViT-Base CLIP encoder as it is closest in size to our ViT-Small encoders. We note that this comparison is imperfect due to difference in data distributions and encoders but believe it is still instructive to see. In Figure 13 we show the results. We observe a promising signal that self-supervised representations can outperform strong CLIP encoders. We believe it would be interesting to perform a controlled study on in-the-wild images with text annotations, like from the recently released Ego4D dataset (Grauman et al., 2021).

MoCo-v3 comparisons. Finally, we compare visual encoders trained with the Momentum Contrastive (MoCo) self-supervised learning framework instead of MAE used in MVP. We opt to use the latest MoCo-v3 (Chen et al., 2021b) designed for ViT models. We show results in Figure 14. We observe that MoCo-v3 can achieve good performance on one of the tasks and non-trivial on the other. This suggests that our approach may be more general and applicable to other self-supervised pre-training techniques as well. However, in contrast to MAE, it may be harder to adapt techniques like MoCo-v3 to in-the-wild images (see Figure 8).

6. Related Work

Dexterous manipulation. Recently, OpenAI (OpenAI et al., 2020; 2019) has shown impressive results in dexterous in-hand manipulation with large scale domain randomization. A number of works consider related manipulation problems with multi-finger hands but rely on explicit state estimation (Handa et al., 2020; Huang et al., 2021), expert policies (Chen et al., 2021a), human demonstrations (Rajeswaran et al., 2018; Radosavovic et al., 2021; Qin et al., 2021), human priors (Mandikal & Grauman, 2021), or models (Nagabandi et al., 2019). In contrast, we do not use on any of the aforementioned components in our approach.

Representations in RL. One way to learn representations for motor control is to rely on the task signal. This is commonly done in end-to-end RL (Mnih et al., 2015; Levine et al., 2016; Kalashnikov et al., 2018). However, it results in high sample complexity, particularly in the case of high-dimensional observations like images. Furthermore, such representations may get overly adapted to the problem at hand and not generalize to new settings (e.g., new objects).

RL with self-supervision. One way to overcome the high sample complexity of RL is to employ auxiliary objectives. In particular, in addition to learning the task, learn to predict some property of the environment, e.g., depth, that may lead to learning good representations as a side effect (Mirowski et al., 2017; Jaderberg et al., 2017; Shelhamer et al., 2017; Lample & Chaplot, 2017). Rather than predicting hand-designed environment properties, researchers explored using more general self-supervised objectives (Oord et al., 2018; Yarats et al., 2019; Srinivas et al., 2020). For example, Srinivas et al. (2020) show excellent performance in vision-based RL tasks. Representations can also be pre-trained on the data from the environment (Ha & Schmidhuber, 2018; Srinivas et al., 2020). Overall, we share the goal of learning good visual representations with self-supervision. In contrast, we learn representations from large collections of natural images rather than environment-specific experience.

Self-supervision in robotics. Self-supervised learning has also been used in various robotic settings. (Pinto & Gupta, 2016) and Agrawal et al. (2016) learn representations through interaction. Sermanet et al. (2018) learn representations from multiview video using contrastive learning and use them for imitation learning. Florence et al. (2018) learn dense image descriptors with self-supervision. Zhan et al. (2020) learn robotic manipulation with RAD (Laskin et al., 2020) and CURL (Srinivas et al., 2020). Pari et al. (2021) show the effectiveness of visual representations with non-parametric nearest neighbor controllers. All of these approaches learn representations in a specific robotic setting of interest (e.g., videos in the lab), rather than general visual representations from image collections like ours.

Supervised pre-training. Sax et al. (2018) and Chen et al. (2020a) show that representations learned from performing a set of mid-level vision tasks using label supervision benefits downstream navigation and manipulation tasks, respectively. Zhou et al. (2019) show the effectiveness on visual representations in driving settings. Yen-Chen et al. (2020) transfer image models trained on supervised vision tasks, e.g., edge detection and semantic segmentation, to affordance prediction models for object manipulation. Shah & Kumar (2021) use ImageNet representations for dexterous manipulation. In contrast to all of these, our approach is self-supervised and does not rely on labeled datasets.

Self-supervised pre-training in computer vision. Self-supervised learning has been gaining momentum in computer vision. The approaches often rely on pretext tasks for pre-training (Doersch et al., 2015; Wang & Gupta, 2015; Noroozi & Favaro, 2016; Zhang et al., 2016; Pathak et al., 2016; Komodakis & Gidaris, 2018). More recently, contrastive learning methods, e.g., (Hadsell et al., 2006; Oord et al., 2018; Wu et al., 2018; Henaff, 2020; He et al., 2020; Chen et al., 2020c; Jabri et al., 2020), have been popular. These techniques try to learn to be invariant to a set of hand-crafted augmentations. Xiao et al. (2021b) have shown that the augmentations introduce inductive bias and may harm downstream transfer. Masked image autoencoding (Chen et al., 2020b; Bao et al., 2022; He et al., 2021) pursues a different direction by learning to recover masked pixels. Specifically, we adopt the Masked Autoencoders (MAE) (He et al., 2021) have shown excellent performance on recognition tasks. We adopt the MAE as our visual pre-training strategy for learning motor control.

7. Conclusion

In this paper, we show that self-supervised visual pre-training is effective for motor control. We use a single vision encoder to learn various motor control tasks from pixels, without per-task fine-tuning, explicit state estimation, or expert demonstrations. We further show large sample complexity improvements compared to supervised baselines (up to 80% absolute success rate) and sometimes even match the oracle state performance. Finally, we show that in-the-wild images, e.g., from YouTube or Egocentric videos, can lead to better visual representations than ImageNet images.

Acknowledgements

We thank William Peebles, Matthew Tancik, Anastasios Angelopoulos, Aravind Srinivas, and Agrim Gupta for helpful discussions. This work was supported in part by DOD including DARPA’s MCS, XAI, LwLL, and/or SemaFor programs; ONR MURI program (N00014-14-1-0671), as well as BAIR’s industrial alliance programs.

References

- Agarwal, R., Schwarzer, M., Castro, P. S., Courville, A. C., and Bellemare, M. Deep reinforcement learning at the edge of the statistical precipice. *NeurIPS*, 2021.
- Agrawal, P., Nair, A. V., Abbeel, P., Malik, J., and Levine, S. Learning to poke by poking: Experiential learning of intuitive physics. *NeurIPS*, 2016.
- Bao, H., Dong, L., and Wei, F. Beit: Bert pre-training of image transformers. In *ICLR*, 2022.
- Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. Language models are few-shot learners. *NeurIPS*, 2020.
- Calli, B., Singh, A., Walsman, A., Srinivasa, S., Abbeel, P., and Dollar, A. M. The ycb object and model set: Towards common benchmarks for manipulation research. In *ICAR*, 2015.
- Chen, B., Sax, A., Lewis, G., Armeni, I., Savarese, S., Zamir, A., Malik, J., and Pinto, L. Robust policies via mid-level visual representations: An experimental study in manipulation and navigation. In *CoRL*, 2020a.
- Chen, M., Radford, A., Child, R., Wu, J., Jun, H., Luan, D., and Sutskever, I. Generative pretraining from pixels. In *ICML*, 2020b.
- Chen, T., Kornblith, S., Norouzi, M., and Hinton, G. A simple framework for contrastive learning of visual representations. In *ICML*, 2020c.
- Chen, T., Xu, J., and Agrawal, P. A system for general in-hand object re-orientation. In *CoRL*, 2021a.
- Chen, X., Xie, S., and He, K. An empirical study of training self-supervised vision transformers. In *ICCV*, 2021b.
- Damen, D., Doughty, H., Farinella, G. M., Fidler, S., Furnari, A., Kazakos, E., Moltisanti, D., Munro, J., Perrett, T., Price, W., et al. Scaling egocentric vision: The epic-kitchens dataset. In *ECCV*, 2018.
- Damen, D., Doughty, H., Farinella, G. M., Furnari, A., Ma, J., Kazakos, E., Moltisanti, D., Munro, J., Perrett, T., Price, W., and Wray, M. Rescaling egocentric vision: Collection, pipeline and challenges for epic-kitchens-100. *IJCV*, 2021.
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HCT*, 2019.
- Doersch, C., Gupta, A., and Efros, A. A. Unsupervised visual representation learning by context prediction. In *ICCV*, 2015.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2020.
- Florence, P. R., Manuelli, L., and Tedrake, R. Dense object nets: Learning dense visual object descriptors by and for robotic manipulation. In *CoRL*, 2018.
- Girshick, R., Donahue, J., Darrell, T., and Malik, J. Rich feature hierarchies for accurate object detection and semantic segmentation. In *CVPR*, 2014.
- Goyal, P., Dollár, P., Girshick, R., Noordhuis, P., Wesolowski, L., Kyrola, A., Tulloch, A., Jia, Y., and He, K. Accurate, large minibatch sgd: Training imagenet in 1 hour. *arXiv:1706.02677*, 2017a.
- Goyal, R., Ebrahimi Kahou, S., Michalski, V., Materzynska, J., Westphal, S., Kim, H., Haenel, V., Freund, I., Yianilos, P., Mueller-Freitag, M., et al. The” something something” video database for learning and evaluating visual common sense. In *ICCV*, 2017b.
- Grauman, K., Westbury, A., Byrne, E., Chavis, Z., Furnari, A., Girdhar, R., Hamburger, J., Jiang, H., Liu, M., Liu, X., et al. Ego4d: Around the world in 3,000 hours of egocentric video. *arXiv:2110.07058*, 2021.
- Ha, D. and Schmidhuber, J. World models. *arXiv:1803.10122*, 2018.
- Hadsell, R., Chopra, S., and LeCun, Y. Dimensionality reduction by learning an invariant mapping. In *CVPR*, 2006.
- Handa, A., Van Wyk, K., Yang, W., Liang, J., Chao, Y.-W., Wan, Q., Birchfield, S., Ratliff, N., and Fox, D. Dexpilot: Vision-based teleoperation of dexterous robotic hand-arm system. In *ICRA*, 2020.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *CVPR*, 2016.
- He, K., Gkioxari, G., Dollár, P., and Girshick, R. Mask r-cnn. In *ICCV*, 2017.
- He, K., Fan, H., Wu, Y., Xie, S., and Girshick, R. Momentum contrast for unsupervised visual representation learning. In *CVPR*, 2020.
- He, K., Chen, X., Xie, S., Li, Y., Dollár, P., and Girshick, R. Masked autoencoders are scalable vision learners. *arXiv:2111.06377*, 2021.

- Henaff, O. Data-efficient image recognition with contrastive predictive coding. In *ICML*, 2020.
- Huang, W., Mordatch, I., Abbeel, P., and Pathak, D. Generalization in dexterous manipulation via geometry-aware multi-task learning. *arXiv:2111.03062*, 2021.
- Jabri, A., Owens, A., and Efros, A. Space-time correspondence as a contrastive random walk. *NeurIPS*, 2020.
- Jaderberg, M., Mnih, V., Czarnecki, W. M., Schaul, T., Leibo, J. Z., Silver, D., and Kavukcuoglu, K. Reinforcement learning with unsupervised auxiliary tasks. In *ICLR*, 2017.
- James, S., Ma, Z., Arrojo, D. R., and Davison, A. J. Rlbench: The robot learning benchmark & learning environment. *RA-L*, 2020.
- Kalashnikov, D., Irpan, A., Pastor, P., Ibarz, J., Herzog, A., Jang, E., Quillen, D., Holly, E., Kalakrishnan, M., Vanhoucke, V., et al. Scalable deep reinforcement learning for vision-based robotic manipulation. In *CoRL*, 2018.
- Klambauer, G., Unterthiner, T., Mayr, A., and Hochreiter, S. Self-normalizing neural networks. In *NeurIPS*, 2017.
- Komodakis, N. and Gidaris, S. Unsupervised representation learning by predicting image rotations. In *ICLR*, 2018.
- Lample, G. and Chaplot, D. S. Playing fps games with deep reinforcement learning. In *AAAI*, 2017.
- Laskin, M., Lee, K., Stooke, A., Pinto, L., Abbeel, P., and Srinivas, A. Reinforcement learning with augmented data. *NeurIPS*, 2020.
- Levine, S., Finn, C., Darrell, T., and Abbeel, P. End-to-end training of deep visuomotor policies. *JMLR*, 2016.
- Makoviychuk, V., Wawrzyniak, L., Guo, Y., Lu, M., Storey, K., Macklin, M., Hoeller, D., Rudin, N., Allshire, A., Handa, A., et al. Isaac gym: High performance gpu-based physics simulation for robot learning. In *NeurIPS*, 2021.
- Mandikal, P. and Grauman, K. Dexvip: Learning dexterous grasping with human hand pose priors from video. In *CoRL*, 2021.
- Mirowski, P., Pascanu, R., Viola, F., Soyer, H., Ballard, A. J., Banino, A., Denil, M., Goroshin, R., Sifre, L., Kavukcuoglu, K., et al. Learning to navigate in complex environments. In *ICLR*, 2017.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., et al. Human-level control through deep reinforcement learning. *Nature*, 2015.
- Nagabandi, A., Konolige, K., Levine, S., and Kumar, V. Deep dynamics models for learning dexterous manipulation. In *CoRL*, 2019.
- Noroozi, M. and Favaro, P. Unsupervised learning of visual representations by solving jigsaw puzzles. In *ECCV*, 2016.
- Oord, A. v. d., Li, Y., and Vinyals, O. Representation learning with contrastive predictive coding. *arXiv:1807.03748*, 2018.
- OpenAI, Akkaya, I., Andrychowicz, M., Chociej, M., Litwin, M., McGrew, B., Petron, A., Paino, A., Plappert, M., Powell, G., Ribas, R., Schneider, J., Tezak, N., Tworek, J., Welinder, P., Weng, L., Yuan, Q., Zaremba, W., and Zhang, L. Solving rubik’s cube with a robot hand. *arXiv:1910.07113*, 2019.
- OpenAI, Andrychowicz, M., Baker, B., Chociej, M., Józefowicz, R., McGrew, B., Pachocki, J., Petron, A., Plappert, M., Powell, G., Ray, A., Schneider, J., Sidor, S., Tobin, J., Welinder, P., Weng, L., and Zaremba, W. Learning dexterous in-hand manipulation. *The International Journal of Robotics Research*, 2020.
- Pari, J., Muhammad, N., Arunachalam, S. P., Pinto, L., et al. The surprising effectiveness of representation learning for visual imitation. *arXiv:2112.01511*, 2021.
- Pathak, D., Krahenbuhl, P., Donahue, J., Darrell, T., and Efros, A. A. Context encoders: Feature learning by inpainting. In *CVPR*, 2016.
- Pinto, L. and Gupta, A. Supersizing self-supervision: Learning to grasp from 50k tries and 700 robot hours. In *ICRA*, 2016.
- Qin, Y., Wu, Y.-H., Liu, S., Jiang, H., Yang, R., Fu, Y., and Wang, X. Dexmv: Imitation learning for dexterous manipulation from human videos. *arXiv:2108.05877*, 2021.
- Radford, A., Narasimhan, K., Salimans, T., and Sutskever, I. Improving language understanding by generative pre-training. 2018.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I., et al. Language models are unsupervised multitask learners. 2019.
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021.
- Radosavovic, I., Wang, X., Pinto, L., and Malik, J. State-only imitation learning for dexterous manipulation. In *IROS*, 2021.

- Rajeswaran, A., Kumar, V., Gupta, A., Vezzani, G., Schulman, J., Todorov, E., and Levine, S. Learning complex dexterous manipulation with deep reinforcement learning and demonstrations. *RSS*, 2018.
- Rudin, N., Hoeller, D., Reist, P., and Hutter, M. Learning to walk in minutes using massively parallel deep reinforcement learning. In *CoRL*, 2021.
- Sax, A., Emi, B., Zamir, A. R., Guibas, L., Savarese, S., and Malik, J. Mid-level visual representations improve generalization and sample efficiency for learning visuomotor policies. *arXiv:1812.11971*, 2018.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms. *arXiv:1707.06347*, 2017.
- Sermanet, P., Lynch, C., Chebotar, Y., Hsu, J., Jang, E., Schaal, S., Levine, S., and Brain, G. Time-contrastive networks: Self-supervised learning from video. In *ICRA*, 2018.
- Shah, R. and Kumar, V. Rrl: Resnet as representation for reinforcement learning. *arXiv:2107.03380*, 2021.
- Shan, D., Geng, J., Shu, M., and Fouhey, D. F. Understanding human hands in contact at internet scale. In *CVPR*, 2020.
- Shelhamer, E., Mahmoudieh, P., Argus, M., and Darrell, T. Loss is its own reward: Self-supervision for reinforcement learning. In *ICLR*, 2017.
- Srinivas, A., Laskin, M., and Abbeel, P. Curl: Contrastive unsupervised representations for reinforcement learning. In *ICML*, 2020.
- Tassa, Y., Doron, Y., Muldal, A., Erez, T., Li, Y., Casas, D. d. L., Budden, D., Abdolmaleki, A., Merel, J., Lefrancq, A., et al. Deepmind control suite. *arXiv:1801.00690*, 2018.
- Todorov, E., Erez, T., and Tassa, Y. Mujoco: A physics engine for model-based control. In *IROS*, 2012.
- Van Den Oord, A., Dieleman, S., Zen, H., Simonyan, K., Vinyals, O., Graves, A., Kalchbrenner, N., Senior, A. W., and Kavukcuoglu, K. Wavenet: A generative model for raw audio. *SSW*, 2016.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. In *NeurIPS*, 2017.
- Wang, Q., Li, B., Xiao, T., Zhu, J., Li, C., Wong, D. F., and Chao, L. S. Learning deep transformer models for machine translation. In *ACL*, 2019.
- Wang, X. and Gupta, A. Unsupervised learning of visual representations using videos. In *ICCV*, 2015.
- Wu, Z., Xiong, Y., Yu, S. X., and Lin, D. Unsupervised feature learning via non-parametric instance discrimination. In *CVPR*, 2018.
- Xiao, T., Dollar, P., Singh, M., Mintun, E., Darrell, T., and Girshick, R. Early convolutions help transformers see better. In *NeurIPS*, 2021a.
- Xiao, T., Wang, X., Efros, A. A., and Darrell, T. What should not be contrastive in contrastive learning. In *ICLR*, 2021b.
- Yarats, D., Zhang, A., Kostrikov, I., Amos, B., Pineau, J., and Fergus, R. Improving sample efficiency in model-free reinforcement learning from images. *arXiv:1910.01741*, 2019.
- Yen-Chen, L., Zeng, A., Song, S., Isola, P., and Lin, T.-Y. Learning to see before learning to act: Visual pre-training for manipulation. In *ICRA*, 2020.
- Yu, T., Quillen, D., He, Z., Julian, R., Hausman, K., Finn, C., and Levine, S. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning. In *CoRL*, 2020.
- Zhan, A., Zhao, P., Pinto, L., Abbeel, P., and Laskin, M. A framework for efficient robotic manipulation. *arXiv:2012.07975*, 2020.
- Zhang, R., Isola, P., and Efros, A. A. Colorful image colorization. In *ECCV*, 2016.
- Zhou, B., Krähenbühl, P., and Koltun, V. Does computer vision matter for action? *Science Robotics*, 2019.
- Zhu, Y., Wong, J., Mandlekar, A., and Martín-Martín, R. robosuite: A modular simulation framework and benchmark for robot learning. *arXiv:2009.12293*, 2020.