

# VideoDex: Learning Dexterity from Internet Videos

Kenneth Shaw\*    Shikhar Bahl\*    Deepak Pathak

Carnegie Mellon University

**Abstract:** To build general robotic agents that can operate in many environments, it is often imperative for the robot to collect experience in the real world. However, this is often not feasible due to safety, time, and hardware restrictions. We thus propose leveraging the next best thing as real-world experience: internet videos of humans using their hands. Visual priors, such as visual features, are often learned from videos, but we believe that more information from videos can be utilized as a stronger prior. We build a learning algorithm, VideoDex, that leverages *visual*, *action*, and *physical* priors from human video datasets to guide robot behavior. These actions and physical priors in the neural network dictate the typical human behavior for a particular robot task. We test our approach on a robot arm and dexterous hand-based system and show strong results on various manipulation tasks, outperforming various state-of-the-art methods. For videos and supplemental material visit our website at <https://video-dex.github.io>

**Keywords:** Dexterous Manipulation, Large Scale Robotics, Imitation Learning

## 1 Introduction

The long-standing dream of many roboticists is to see robots autonomously perform diverse tasks in diverse environments. To build a robot that can operate anywhere, many methods rely on successful robotic interaction data to train on. However, deploying inexperienced, real-world robots to collect experience may require constant supervision which is infeasible. This poses a chicken-and-egg problem for robot learning because to collect experience safely, the robot already needs to be experienced. How do we get around this deadlock?

Fortunately, there is plenty of real-world human interaction videos on the internet. This data can potentially help bootstrap robot learning by side-stepping the data collection-training loop. This insight of leveraging human videos to aid robotics is not new and has seen immense attention from the community at large [1, 2, 3]. However, most of the prior work tends to use human data as a mechanism for pretraining just the visual representation [4, 5, 6, 7, 8], much like how deep learning has been used as a pretraining tool in related areas of computer vision [9, 10] and natural language processing [11, 12]. Although pretraining visual representations can aid in efficiency, we believe that a large part of the inefficiency stems from very large action spaces. For continuous control, learning this is exponential in the number of actions and timesteps, and even more difficult for high degree-of-freedom robots (shown in Figure 1). Dexterous hands are one such class of high degree of freedom robots that have the possibility to provide great contact for the grasping and manipulation of different objects. Their similarity to human hands makes learning from human video advantageous.

In this work, we study how to go beyond using internet human videos merely as a source of visual pretraining (i.e. **visual priors**), and leverage the information of how humans move their limbs to guide train robots on how they should move (i.e. **action priors**). However, guiding robot motions using human videos requires understanding the scene in 3D, figuring out human intent, and transferring from human to robot embodiment. First, 3D human estimation works decently well in general human videos which we can leverage to gather 3D understanding. Second, there have been large-scale datasets that break down the human intent via crowdsourcing labels [2, 1]. Finally, to handle the embodiment transfer, we use human hand to robot hand retargeting as an energy function to pretrain the robot action policy. Our key insight is to combine these visual and action priors from human videos with a prior on how robot should move in the world [13, 14] (i.e., **physical prior**, using a second order dynamical system) to obtain dexterous robot policies that can act in the real world. We

---

\*Equal contribution, order decided by coin flip.

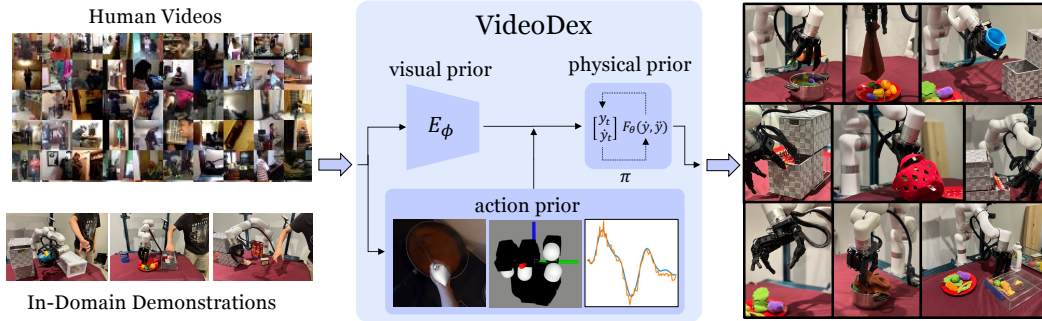


Figure 1: We re-target human videos as an action prior, use pretrained embeddings as a visual prior, and use Neural Dynamical Policies (NDPs) [13] as a physical prior to complete many different tasks on a robotic hand.

call this approach, VideoDex. To enhance real-world performance, we mix the experience obtained from massive internet data with a few in-domain demonstrations.

In summary, VideoDex is a robot learning algorithm that incorporates visual, action, and physical priors into a single open-loop policy by learning from passive videos contained in human activity datasets from the internet. VideoDex then only needs to adapt to real world tasks using a few in-domain examples. We find that VideoDex outperforms many state-of-the-art robot learning methods on seven different real-world manipulation tasks on a high DOF multi-fingered robotic arm-hand system as well as on a 1-DOF gripper robotic arm system.

## 2 Related Work

**Learning for Dexterity** Reinforcement learning (RL) with an engineered reward function can show dexterous simulation results [15, 16] but requires lots of data, especially in high DOF dexterous manipulation. This requires simulators [17, 18], which cannot model physics properly, making real-world transfer difficult. Behavior cloning is an approach [19, 20] that can work safely. DIME [21] involves using nearest neighbor matching of image representations with demonstrations to determine actions. Qin et al. [22] teleoperates and learns policies in simulation, followed by Sim2Real transfer. DexMV[23] uses collected human hand videos for robot hand imitation learning. DexVIP [24] learns hand-object affordances and priors for RL initialization using curated video datasets.

**Learning from Videos and Large-Scale Datasets** There are many curated datasets from internet human videos, for example, FreiHand [25] for hand poses, 100 Days of Hands [26] for hand-object interactions, Something-Something [3] for semantically similar interactions, Human3.6M [27] and the CMU Mocap Database [28] for Human pose estimation. Epic Kitchens [2], ActivityNet datasets [29], or YouCook [30] are action-driven datasets we focus on for dexterous manipulation.

**Learning Action from Videos** Detecting humans, estimating poses of different body parts, or understanding the dynamics and interactions related to human motion is a commonly studied problem. One can model human hands using the MANO [31] model and the human body using SMPL, SMPL-X [32, 33] models. There are many efforts in human pose estimation such as [34, 35, 36]. We focus on FrankMocap [36] for our project as it is robust for online videos. Traditionally, teleoperation approaches have employed hand markers with gloves for motion capture [37] or VR settings [38]. Without gloves, Li et. al. [39] used depth images and a paired human-robot dataset for teleoperation, and Handa et. al. [40] designed a system that mimics the functional intent of the human operator to perform object manipulation tasks.

**Robot Learning by Watching Humans** Recent works have leveraged human datasets to learn cost functions [41, 42, 43], learn action correspondences [44] both in a paired [45] and unpaired manner [46]. This data can also be used to extract explicit actions by leveraging structure in the collection (such as reacher-grabber tools [47]) or prediction of future hand and object locations [48], as well as keypoint detectors [49]. This can also be used to build representations for robot learning [6, 50]. R3M [6] trains on the Ego4D [1] dataset using a temporal alignment loss between language labels and video frames. We build on top of previous efforts in this area, where we combine visual representations trained on human activity data, with *action* driven representations.



Figure 2: The collection of train objects (left) and test objects (right) used for experimentation.

### 3 Background

#### 3.1 Neural Dynamic Policies

Neural Dynamic Policies (NDPs) [13, 14, 51], produce smooth and safe open-loop trajectories. When using them as a network backbone, they can be rolled out to trajectories of arbitrary lengths which enables the use of varying-length human videos. NDPs can be described with the Dynamic Movement Primitive equation [52, 53, 54, 55]:

$$\ddot{y} = \alpha(\beta(g - y) - \dot{y}) + f_w(x, g), \quad (1)$$

where  $y$  is the coordinate frame of the robot,  $g$  is the desired goal in the given coordinate frame,  $f_w$  is a radial basis forcing function,  $x$  is a time variable, and  $\alpha, \beta$  are global constants. NDPs use the robot state, scene, and a NN to output the goal  $g$  and shape parameters  $w$  of the forcing function  $f_w$ .

#### 3.2 Learning from Watching Humans

Recently, Sivakumar et al. [56] introduced Robotic Telekinesis, a pipeline that teleoperates the Allegro Hand [57] using a single RGB camera. Leveraging work in monocular human hand and body pose estimation [36], hand and body modeling [31, 32, 33], and human internet data, Robotic Telekinesis real-time re-targets the human hand and body to the robot hand and arm. Due to its efficiency and ease of use, we leverage Sivakumar et al. [56]’s approach for demonstration collection.

We borrow the human hand to robot hand re-targeting method from Robotic Telekinesis [56] that manually defines key vectors  $v_i^h$  and  $v_i^r$  between palms and fingertips on both the human and robot hand. They build an energy function  $E_\pi$  which minimizes the distance between human hand poses  $(\beta, \theta)$  and robot hand poses  $q$ .  $c_i$  is a scale parameter. Therefore, the energy function is defined as:

$$E_\pi((\beta_h, \theta_h), q) = \sum_{i=1}^{10} \|v_i^h - (c_i \cdot v_i^r)\|_2^2 \quad (2)$$

Sivakumar et al. [56] train an MLP  $H_R(\cdot)$  to implicitly minimize this energy function in 2, conditioned on knowing human poses  $(\beta, \theta)$ . For more details, we refer the readers to Sivakumar et al. [56].

### 4 Learning Dexterity from Human Videos

We learn general-purpose manipulation by utilizing large-scale human hand action data as prior robot experience. We leverage not only visual priors of the scene’s appearance but also leverage important aspects of the human hand’s motion, intent, and interaction. To do this, we *re-target* the human video data to trajectories from the robot’s embodiment and point of view. By pretraining policies with these human hand trajectories, we learn *action* priors on how the robot should behave. However, it’s notoriously difficult to leverage these noisy human video detections. Therefore, we must also employ a policy with *physical* priors to learn smooth and robust policies that do not overfit to noise. We explain insights and our method used to leverage *action* priors in the sections below.

#### 4.1 Visual Priors from Human Activity Data

Many previous works [6, 7, 8] have tackled visual priors and representations for robot learning. These networks often encode some form of semantic visual priors into the pretrained network from human video internet datasets. We use the encoder from Nair et al. [6] as a useful visual initialization for our policy. Nair et al. [6] is trained on a visual-language alignment as well as a temporal consistency loss. Our network takes human video frames and processes them using the publicly released ResNet18 [58] encoder,  $E_\phi$  from R3M [6]. The output of this network is our visual representation for learning.

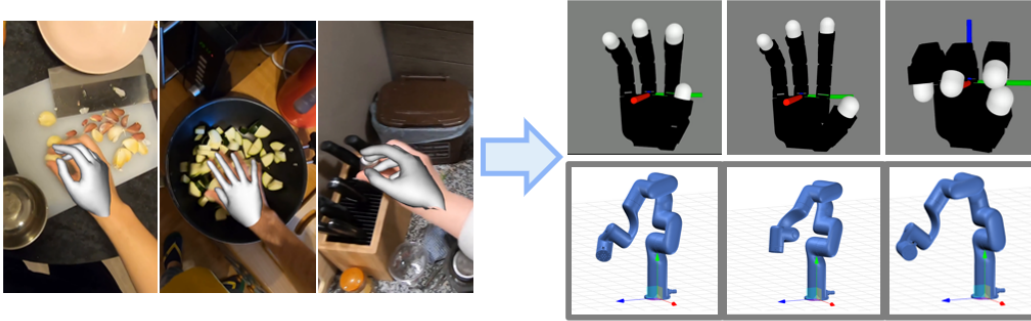


Figure 3: To use internet videos as pseudo-robot experience, we re-target human hand detections from the 3D MANO model [31] to 16 DoF robotic hand (LEAP) embodiment and we re-target the wrist from the moving camera to the xArm6 [59] embodiment. Videos at <https://video-dex.github.io>

## 4.2 Action Priors from Human Activity Data

While visual pretraining aids in semantic understanding, human data contains a lot more information about how to interact with the world. VideoDex uses action information to pretrain an action prior, a network initialization that encodes information about the typical actions for a particular task.

However, training robot policies on human actions are difficult, as there is a large embodiment gap between humans and robots as described in Handa et al. [40] and Sivakumar et al. [56] Thus, we must re-target the motion of the human to the robot embodiment to use it in training. This problem is solved using three main components. First, we detect human hands in videos. Second, we project hand poses  $H$  to robot finger joints  $H_r$ . Finally, we convert human wrist pose  $P$  to robot arm pose  $P_r$ .  $H_r$  and  $P_r$  define the trajectory of the human in the robot’s frame, from which we can extract actions to pretrain our policy network with the action prior. See Figure 4 for a summary of the stages.

**Action and Hand Detections** First, we must detect the right actions the human is completing. To expedite development, we use the action annotations from the EpicKitchens dataset [2] but an action detection network such as [60] can be used. Now, we must detect the hand. VideoDex first computes a crop  $c$  around the operator’s hand using OpenPose [61] and the result is passed to FrankMocap [36] to obtain hand shape ( $\beta$ ) and pose parameters ( $\theta$ ) of the 3D MANO model [31]. These parameters are passed through a low pass filter and subsequently used in re-targeting to the robot.

**Re-targeting Wrist Pose** In this section, we show how to compute the transformation that describes the wrist pose in the robot frame denoted as  $M_{Robot}^{Wrist}$ . First, to calculate  $M_{C_t}^{Wrist}$ , where  $C_t$  is the camera frame at timestep  $t$  we leverage the Perspective-n-point algorithm [62]. This takes 2D keypoint outputs  $(u_i, v_i)$  by the hand detection model and 3D keypoints from the hand model  $(x_i, y_i, z_i)$  and computes  $M_{C_t}^{Wrist}$ . To accurately obtain camera intrinsics for PnP, COLMAP is used [63].

In human egocentric video datasets, the position of the camera is not fixed and we must compensate for this movement. Specifically, we compute the transformation between the camera pose in the first frame  $C_1$  and all other frames in the trajectory,  $C_t$ . We call this transform  $M_{C_1}^{C_t}$ . To estimate this, we run monocular SLAM, specifically ORBSLAM3 [64].

Computing wrist poses in the first camera coordinate frame is important but this is still not in the robot frame because the robot is always upright. To be able to transform the human trajectory in the robot’s frame, we must find the vector that is parallel to gravity in the camera’s frame,  $\alpha_p$ . Thus recover object segmentations for surfaces that are parallel to the floor such as tables, floors, counters, and similar synonyms using a state-of-the-art object detector (Detic [65]). Then an estimated depth map from RGB frames only using Adabins [66] is computed. This way, the method does not rely on the long-term contiguity of a video like most SLAM approaches. We then use depth map portions that correspond to the relevant objects and calculate a surface normal vector. We estimate  $\alpha_p$  using this normal vector and the following equations:

$$\text{pitch} = \tan^{-1}(x_{Acc}/\sqrt{y_{Acc}^2 + z_{Acc}^2}) \quad (3)$$

$$\text{roll} = \tan^{-1}(y_{Acc}/\sqrt{x_{Acc}^2 + z_{Acc}^2}) \quad (4)$$

Detailed ablations on the parameterization of the initial pitch of the predicted trajectory ( $\alpha$ ) are provided in Section 6. In SLAM, we also remove the dependency on gyroscope data by assuming

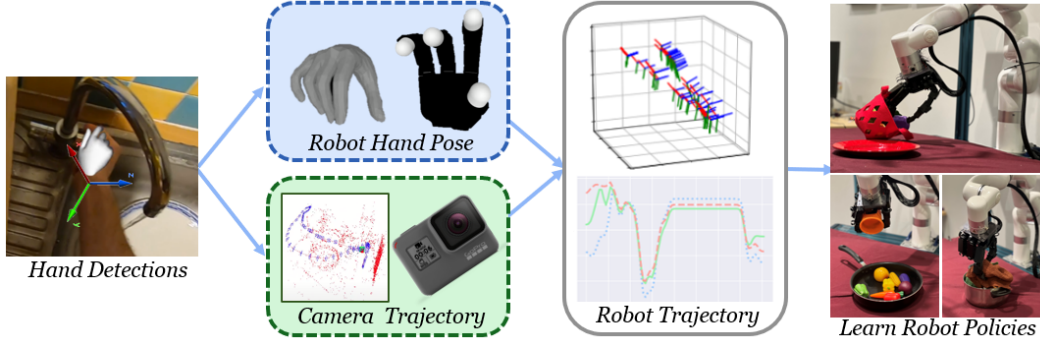


Figure 4: To use human videos as an action prior for training policies, we re-target them to the robot embodiment. The detected human fingers are converted to the robot fingers using a learned energy function. The wrist is re-targeted using the detections and camera trajectory and transformed to the robot arm.

that the scaling factor is 1.0. This is acceptable because the trajectory is rescaled to the robot frame later. Therefore, this wrist re-targeting approach uses only 2D images from human videos.

Since the robot has workspace limits, and we would also like to center the starting pose of the robot, we heuristically compute  $T_{Robot}^{World}$  which rescales and rotates the human trajectory in the world frame  $\tau_W^{wrist}$  into the robot trajectory  $\tau_R^{wrist}$ . The final function to obtain  $M_{Robot}^{Wrist}$  can be described as:

$$M_{Robot}^{Wrist} = T_{Robot}^{World} \cdot M_{World}^{C_1} \cdot M_{C_1}^{C_t} \cdot M_{C_t}^{wrist} \quad (5)$$

---

#### Algorithm 1 Procedure for VideoDex

---

**Require:** Human videos  $V_{1:K}^H$  (length  $T$ ), policy  $\pi_\theta$ , demonstrations  $\mathcal{D}_{1:N}$ . Human detection  $f_{human}$  [36].

**for**  $k = 1 \dots K$  **do**

**for**  $t = 1 \dots T$  **do**

    Pose parameters  $\theta_t, \beta_t = f_{human}(I_t)$

    Get wrist pose  $w_t$  from 3, 4 and 5,

    Hand pose  $h_t = H(\theta_t, \beta_t)$

**end for**

  Store all  $h_t, w_t$  into robot trajectory  $\tau_R^k$

$\hat{\tau}_R^k = \pi_\theta(I_1^k, h_1^k, w_1^k)$

  Optimize  $\mathcal{L}_\theta = \|\tau_R^k - \hat{\tau}_R^k\|_1$

**end for**

Store policy weights  $\theta_h$  to initialize  $\pi_\theta$

**while** not converged **do**

**for**  $n = 1 \dots N$  **do**

$\tau_n, I_{1:T}^n = \mathcal{D}_n$

$\hat{\tau}_n = \pi_\theta(I_1^n, h_1^n, w_1^n)$

    Optimize  $\mathcal{L}_\theta = \|\tau_n - \hat{\tau}_n\|_1$

**end for**

**end while**

---

**Re-targeting Hand Pose** Human hands are also in a different *embodiment* compared to that of robot hands, like our 16 DOF LEAP Hand. Similarly, to Sivakumar et al. [56], we use  $H(\cdot)$  to map hand poses to robot hand poses. Given human detected pose  $x_h$ , we obtain  $x_r = H(x_h)$  using a similar re-targeting network to Sivakumar et al. [56], and get human hand trajectories:  $\tau_R^{hand}$  in the robot’s embodiment. We use  $\tau_R$  to denote the combined hand and wrist trajectories:  $\tau_R^{hand}, \tau_R^{wrist}$ . See Figure 3 for a visualization.

### 4.3 Learning with Human Videos

We must design an open-loop policy  $\pi$  that learns first from the re-targeted human trajectories (the action prior) and then from real robot trajectories collected in teleoperation. Naively, training a neural network policy on  $\tau_R$  will lead to overfitting to noisy hand detections. To circumvent this, we first use visual priors from the visual ResNet-based [58] encoder provided by Nair et al. [6],  $E_\phi$ . Then, we introduce a *physical prior* to the network, the physically-inspired Neural Dynamic Policies [13, 14].

We construct  $\pi$  with the following setup. We first process the first scene image  $I$  with the visual encoder  $E_\phi$ . Then the extracted features  $E_\phi(I)$  are used to condition an NDP for the wrist and hand separately,  $f_{wrist}$  and  $f_{hand}$ . Concretely, each NDP operates by processing the input features with a small MLP which outputs  $w, g$  which are the trajectory shape and goal parameters. The forward integrator of the NDP outputs an open-loop trajectory for the hand and the wrist,  $\hat{\tau}_R$ . We use the following loss function:

$$\mathcal{L} = \sum_k \text{Loss}_{L1}(\tau_R - [f_{hand}(E_\phi(I_k)), f_{wrist}(E_\phi(I_k))])$$

**Training Methodology:** We use between 500-3000 video clips of humans completing the same task category as the robot will from the Epic Kitchens dataset [2]. For example, in pick, there are

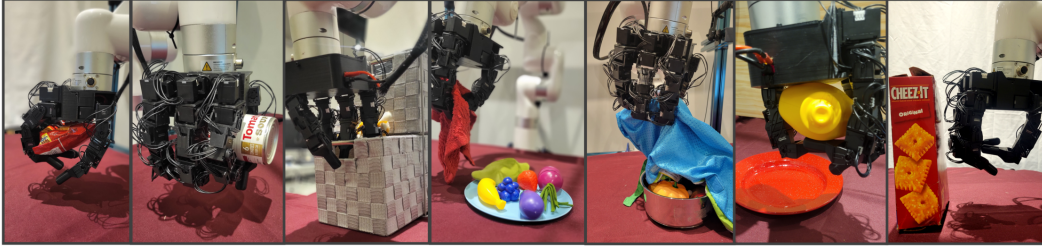


Figure 5: Tasks used in experiments. From left to right: pick, rotate, open, cover, uncover, place and push. See <https://video-dex.github.io> for videos of these tasks.

close to 3000 video clips of humans picking items. These are retargeted to the robot domain and used to pretrain the network with the human action prior of the pick task. Then, the final policy  $\pi$  is trained on a few teleoperated demonstrations of pick on the real robot. The full training takes about 10 hours on a single 2080Ti GPU. More training details can be found in the appendix and in Algorithm 1. Our network consists of the R3M [6] initialized ResNet-18 [58]. We process these features with a 3-layer MLP with a hidden layer size of 512, which are then processed by 2 NDP [13] networks.

## 5 Experimental Setup

We perform thorough real world experiments on manipulation tasks, specifically many tasks that require dexterity. See [our webpage](#) for result videos. We aim to answer the following questions. (1) Is VideoDex able to perform general purpose open-loop manipulation? (2) How much does the action prior of VideoDex help? (3) How much does the physical prior of the NDPs in VideoDex help? (4) What important design choices are there (visual priors, physical priors, or training setup)?

**Task Setup** We pretrain action priors on retargeted Epic Kitchens data for seven robot tasks. Then, we collect about 120-175 demonstrations for each of these tasks on our setup to train the policy. In `pick`, the goal is to pickup an object. In `rotate`, the agent grasps and rotates the object in place. In `cover` and `uncover`, the goal is to cover or uncover a pan/plate with a soft cloth object. `Push` involves flicking/poking an object with the fingers. In `place`, the robot has to pick up an object and place it into a plate, pan or pot. In `open` we open three different drawers. Our testing procedure consists of unseen locations and objects. Details on the tasks and objects are in the supplemental.

While robot hands can provide great dexterity, we also investigate whether 2-finger grippers can benefit from action priors. The internet data is converted to where the closed human hand is a closed 2-finger gripper, and the open human hand is an open 2-finger gripper. We collect separate demonstrations on the real-robot using the 2-finger gripper from xArm [59]. Separate action priors are trained for the 16 DoF LEAP Hand and the 2-finger gripper.

## 6 Results

First, we evaluate the need for initialization with the action priors obtained from the human internet videos.  $\theta_h$ . The baseline without internet pre-training is called BC-NDP. It uses the same physical prior and visual network initialization, without the initialization from  $\theta_h$ . We also compare the effect of the action prior on 2-finger gripper policies. Second, we compare against two standard open-loop behavior cloning approaches introduced in recent benchmarks [51]. BC-open uses a 2 layer MLP instead of the NDP network. BC-RNN, uses an RNN to pre-process the visual features and then a two-stream, 2 layer MLP for wrist and hand trajectories. We try an offline RL ablation CQL [67], where we

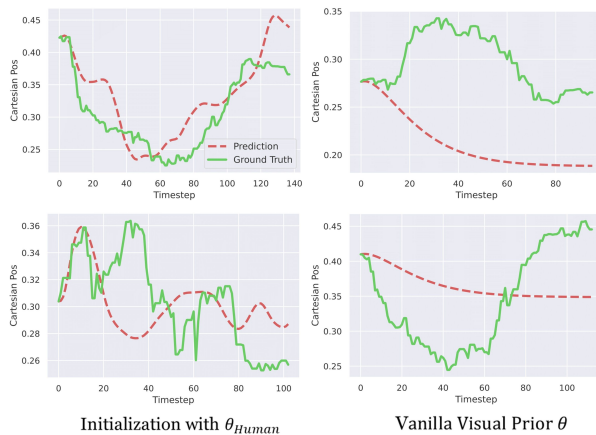


Figure 6: Networks initialized using action priors on human data without further training are closer to ground truth robot trajectories than networks only initialized using visual priors.

	Pick		Rotate		Open		Cover		Uncover		Place		Push	
	train	test	train	test	train	test	train	test	train	test	train	test	train	test
BC-NDP [14]	0.64	0.38	<b>0.94</b>	0.56	<b>0.90</b>	0.60	<b>0.78</b>	0.58	0.88	0.82	0.70	0.35	1.00	0.71
BC-Open[51]	0.50	0.44	0.72	0.38	0.80	0.40	0.44	0.58	<b>1.00</b>	<b>0.91</b>	0.40	0.25	1.00	0.93
BC-RNN [51]	0.56	0.31	0.78	0.50	<b>0.90</b>	0.50	0.56	0.42	0.88	0.75	0.70	0.50	1.00	<b>1.00</b>
VideoDex	<b>0.83</b>	<b>0.77</b>	0.85	<b>0.71</b>	<b>0.80</b>	<b>0.80</b>	<b>0.75</b>	<b>0.63</b>	<b>0.96</b>	0.92	<b>0.89</b>	<b>0.80</b>	1.00	<b>1.00</b>

Table 1: We present the results of train objects and test objects for Videodex and baselines as described above.

use the demonstrations as a sparse reward. We train a behavior cloning policy with the action prior from human videos without the physical prior of the NDP. We call this VideoDex-BC-Open. We ablate the type of visual representation and prior use by trying an initialization using the VGG16 network [68] (VideoDex-VGG) and the MVP network [7] [69] (VideoDex-MVP) based representation trained for robot learning. We ablate the need for a two stream policy, instead training a single NDP for both hand and wrist. (VideoDex-Single) To see if VideoDex works with fewer demonstrations (around 50 demonstrations, 5-7 per variant only), we train a policy called VideoDex-Constrained.

We analyze the results of our experiments and the guiding questions discussed in Section 5. We present the results of our findings as a 0-1 success rate in Table 1 and the result of the ablations we ran on the place task in Table 4.

**Effect of Action Priors** We firstly compare VideoDex against methods that do not employ an action prior trained on human data, as explained in Section 5. For almost all of the tasks, VideoDex either outperforms baselines or has a similar performance, especially for held out objects/instances. We believe that one of the key aspects of VideoDex generalizing to test objects is the action prior pretraining on human videos. This can be seen in Figure 6. Without ever training on the robot demonstrations, the trajectories initialized using the action prior pretrained network  $\theta_h$  (left) are much closer to the ground truth trajectories of a network that is initialized using only a visual prior such as the encoder from Nair et al. [6] (right). From the results, we see that VideoDex-BC-Open with action priors (Table 4) outperforms BC-Open. Having a physical prior added (BC-NDP) tends to help, but it is not the case for every task. We suspect that some tasks require smoother behavior than others. Additionally, in Table 4 our offline RL baseline, CQL [67] does not perform as well as the rest of the approaches, even under-performing the Behavior Cloning setup. Qualitatively, we see a much less smooth and less safe execution with this method, thus we only perform it on one task (place). Note that we use the same visual prior for this as well.

		Place	Open	Pick
1-DOF	BC-Open[51]	0.62	0.69	0.71
1-DOF	VideoDex	<b>0.69</b>	<b>0.82</b>	<b>0.77</b>

Table 2: We compare how the 1-DOF xArm gripper performs using Videodex. [59] Separate demonstrations were collected using this gripper.

	Place	Cover	Uncover
VideoDex-Fixed	0.55	0.50	0.77
VideoDex-Random	0.45	0.63	0.85
VideoDex-IMU	0.70	<b>0.67</b>	0.90
VideoDex	<b>0.80</b>	0.63	<b>0.92</b>

Table 3: Ablations that compare the different ways of calculating the initial pitch of the camera with respect to gravity, on test objects. This enables us to transform human trajectories to be upright like the robot is.

angle. VideoDex-IMU uses the internal image stabilization sensor data to estimate the upright vector. None of these approaches use gyroscope data in SLAM, as we assume that the scaling factor is 1.0. In Table 3, we present the results of these experiments. The performance degrades when randomizing or setting  $M_{World}^{C_1}$  to a fixed value, in all three of the tasks, but it is still comparable to or better than our baselines that do not use any human action data. A possible explanation for the fact that

**Hand vs 2-Finger Gripper** We compare whether the action priors from VideoDex also help in the more general 1-DOF gripper setting. In Table 2, we find that in the 1-DOF setting, VideoDex still improves performance on these tasks. This is because the priors from human internet videos still encode typical wrist trajectory behaviors as well as when the gripper should close for each task.

**Initial Pose Computation Comparison** We compare three different ways to estimate  $\alpha_p$  or  $M_{World}^{C_1}$ , the vector that points parallel to gravity. These methods contrast with VideoDex which uses the surface normal of objects that are typically parallel with the floor to calculate the direction of gravity. VideoDex-Fixed, assumes that  $\alpha_p$  is [0,0]. This is reasonable as we are not relying on robots to exactly mimic the human but get a general action prior. VideoDex-Random, randomizes  $\alpha_p$  in the range of 15-45 degrees, which is the typical egocentric camera

VideoDex-Surface performed better than our VideoDex-IMU is that the sensor data may be noisy and estimating surface normals from visual features is more robust.

**Effect of Physical Priors and Architectural Choices** We compare different types of physical priors in Table 1 and in Table 4. In general (BC-NDP) tends to outperform baselines without a physical prior, except for BC-RNN in a couple of tasks. BC-RNN performs less aggressive behavior, which allowed it to efficiently grasp more objects. In Table 4 it’s shown that an important physical prior is to treat the wrist and the hand in a more disentangled manner, as the performance for VideoDex-Single tends to drop compared to BC-NDP and VideoDex-BC-Open (Behavior Cloning with our action prior pretraining). The two stream architecture aids in learning, as it allows the policy to disentangle the actions of the wrist and the hand. This is important as the same grasp might be used for picking objects in many different locations, and similarly, it is possible to localize many objects and perform completely different types of interactions.

**Generalization with Less Data** We limit VideoDex to a maximum of 5 and 10 teleoperated demonstrations per variant (we have 12-15 variants in our setup). As shown in Table 1, even with 5 instances per variant, we still see a 30% success rate for unseen objects. Empirically, the policies generally go to the right area but are not able to grasp objects properly. With less robot experience, VideoDex outperforms which demonstrates that action priors also boosts sample efficiency.

**Effect of Visual Priors** We compared using our approach with MVP (VideoDex-MVP) [7] and VGG (VideoDex-VGG) [68] and their performance was below VideoDex using Nair et al. [6]. This is likely because both encoders are much larger than the ResNet18 [58] we use and require a lot more training time than feasible on human videos. However, VideoDex-MVP still performs better than VideoDex-VGG, which indicates that using a visual prior trained on human data does in fact help, as Xiao et al. [7] trained the representation in self-supervised fashion on videos and use the embeddings to perform robotics tasks in simulation. We see in Table 1, that while visual priors are important, action priors are more impactful.

**Choice of Robotic Hand** In our experiments, we also tried using the Allegro Hand [57]. We found that the Allegro had higher inaccuracy in control and more hardware failures as compared to LEAP Hand. LEAP Hand outperformed the Allegro Hand 7 – 12% on average in all experiments, thus we use it for our setup.

## 7 Discussion and Limitations

Although we see strong results on the held-out objects, VideoDex has several limitations and scope for future work. First, we focus on curated human video datasets, such as EpicKitchens [2], but only use these as a convenience to expedite our process. It is possible to filter internet videos of humans according to tasks using action detectors and then processing them with VideoDex. We also use camera data in VideoDex but show that with a heuristic driven approach it is possible to obtain similar or better results. Second, we rely on off-the-shelf human hand detection modules that very often have erroneous 6D pose detections, especially when the hand is interacting with objects. Second, the action priors rely on the arm trajectory as well as the hand trajectory retargeting which must be recomputed for each different set of robot parameters and embodiment. Finally, our method of behavior cloning in the real world is currently open-loop, so it cannot react to changes in the environment. This is because closed-loop behavior cloning is difficult to keep safe in the real world. Similarly, when running closed-loop RL it is difficult to guarantee the safety of the system. We leave this to future work, to train policies that can react to changes in the real world.

### Acknowledgments

We thank Aditya Kannan and Shivam Duggal for assisting in robot data collection. We thank Aravind Sivakumar, Russell Mendonca, Jianren Wang, and Sudeep Dasari for fruitful discussions. KS is supported by NSF Graduate Research Fellowship under Grant No. DGE2140739. The work is supported by Samsung GRO Research Award and ONR N00014-22-1-2096.

	Train	Test
<i>Baselines:</i>		
BC-NDP [14]	0.70	0.35
BC-Open [51]	0.40	0.25
BC-RNN [51]	0.70	0.50
CQL [67]	0.40	0.20
<i>No Physical Prior:</i>		
VideoDex-BC-Open	0.50	0.50
VideoDex-Single	0.50	0.30
<i>Visual Prior Ablation:</i>		
VideoDex-VGG	0.20	0.20
VideoDex-MVP	0.40	0.20
<i>Constrained Data:</i>		
VideoDex-Const-5	0.80	0.60
VideoDex-Const-10	0.50	0.30
VideoDex (ours)	<b>0.90</b>	<b>0.70</b>

Table 4: We present the results of the ablations discussed in Section 5. These are all performed on the place task.



## References

- [1] K. Grauman, A. Westbury, E. Byrne, Z. Chavis, A. Furnari, R. Girdhar, J. Hamburger, H. Jiang, M. Liu, X. Liu, et al. Ego4d: Around the world in 3,000 hours of egocentric video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18995–19012, 2022.
- [2] D. Damen, H. Dougherty, G. M. Farinella, S. Fidler, A. Furnari, E. Kazakos, D. Moltisanti, J. Munro, T. Perrett, W. Price, and M. Wray. Scaling egocentric vision: The epic-kitchens dataset. In *European Conference on Computer Vision (ECCV)*, 2018.
- [3] R. Goyal, S. Ebrahimi Kahou, V. Michalski, J. Materzynska, S. Westphal, H. Kim, V. Haenel, I. Freund, P. Yianilos, M. Mueller-Freitag, F. Hoppe, C. Thureau, I. Bax, and R. Memisevic. The “something something” video database for learning and evaluating visual common sense. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- [4] L. Pinto, D. Gandhi, Y. Han, Y.-L. Park, and A. Gupta. The curious robot: Learning visual representations via physical interactions. In *ECCV*, 2016.
- [5] P. Sermanet, C. Lynch, Y. Chebotar, J. Hsu, E. Jang, S. Schaal, and S. Levine. Time-contrastive networks: Self-supervised learning from video. In *ICRA*, 2018.
- [6] S. Nair, A. Rajeswaran, V. Kumar, C. Finn, and A. Gupta. R3m: A universal visual representation for robot manipulation. *arXiv preprint arXiv:2203.12601*, 2022.
- [7] T. Xiao, I. Radosavovic, T. Darrell, and J. Malik. Masked visual pre-training for motor control. *arXiv preprint arXiv:2203.06173*, 2022.
- [8] A. V. Nair, V. Pong, M. Dalal, S. Bahl, S. Lin, and S. Levine. Visual reinforcement learning with imagined goals. In *NeurIPS*, pages 9191–9200, 2018.
- [9] K. He, G. Gkioxari, P. Dollar, and R. Girshick. Mask r-cnn. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- [10] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton. A simple framework for contrastive learning of visual representations. In H. D. III and A. Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 1597–1607. PMLR, 13–18 Jul 2020. URL <https://proceedings.mlr.press/v119/chen20j.html>.
- [11] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei. Language models are few-shot learners. 2020.
- [12] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [13] S. Bahl, M. Mukadam, A. Gupta, and D. Pathak. Neural dynamic policies for end-to-end sensorimotor learning. In *NeurIPS*, 2020.
- [14] S. Bahl, A. Gupta, and D. Pathak. Hierarchical neural dynamic policies. *RSS*, 2021.
- [15] D. Kalashnikov, A. Irpan, P. Pastor, J. Ibarz, A. Herzog, E. Jang, D. Quillen, E. Holly, M. Kalakrishnan, V. Vanhoucke, et al. Qt-opt: Scalable deep reinforcement learning for vision-based robotic manipulation. *arXiv preprint arXiv:1806.10293*, 2018.
- [16] S. Levine, C. Finn, T. Darrell, and P. Abbeel. End-to-end training of deep visuomotor policies. *JMLR*, 2016.
- [17] E. Todorov, T. Erez, and Y. Tassa. MuJoCo: A physics engine for model-based control. In *IROS*, 2012.

- [18] V. Makoviychuk, L. Wawrzyniak, Y. Guo, M. Lu, K. Storey, M. Macklin, D. Hoeller, N. Rudin, A. Allshire, A. Handa, et al. Isaac gym: High performance gpu-based physics simulation for robot learning. *arXiv preprint arXiv:2108.10470*, 2021.
- [19] D. A. Pomerleau. Alvin: An autonomous land vehicle in a neural network. In D. Touretzky, editor, *Advances in Neural Information Processing Systems*, volume 1. Morgan-Kaufmann, 1988. URL <https://proceedings.neurips.cc/paper/1988/file/812b4ba287f5ee0bc9d43bbf5bbe87fb-Paper.pdf>.
- [20] M. Bojarski, D. Del Testa, D. Dworakowski, B. Firner, B. Flepp, P. Goyal, L. D. Jackel, M. Monfort, U. Muller, J. Zhang, X. Zhang, J. Zhao, and K. Zieba. End to end learning for self-driving cars, 2016. URL <https://arxiv.org/abs/1604.07316>.
- [21] S. P. Arunachalam, S. Silwal, B. Evans, and L. Pinto. Dexterous imitation made easy: A learning-based framework for efficient dexterous manipulation, 2022. URL <https://arxiv.org/abs/2203.13251>.
- [22] Y. Qin, H. Su, and X. Wang. From one hand to multiple hands: Imitation learning for dexterous manipulation from single-camera teleoperation, 2022. URL <https://arxiv.org/abs/2204.12490>.
- [23] Y. Qin, Y.-H. Wu, S. Liu, H. Jiang, R. Yang, Y. Fu, and X. Wang. Dexmv: Imitation learning for dexterous manipulation from human videos. *arXiv preprint arXiv:2108.05877*, 2021.
- [24] P. Mandikal and K. Grauman. Dexvip: Learning dexterous grasping with human hand pose priors from video. In *Conference on Robot Learning*, pages 651–661. PMLR, 2022.
- [25] C. Zimmermann, D. Ceylan, J. Yang, B. Russell, M. Argus, and T. Brox. Freihand: A dataset for markerless capture of hand pose and shape from single rgb images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 813–822, 2019.
- [26] D. Shan, J. Geng, M. Shu, and D. F. Fouhey. Understanding human hands in contact at internet scale. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9869–9878, 2020.
- [27] C. Ionescu, D. Papava, V. Olaru, and C. Sminchisescu. Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *IEEE transactions on pattern analysis and machine intelligence*, 36(7):1325–1339, 2013.
- [28] Cmu graphics lab motion capture database. <http://mocap.cs.cmu.edu/>.
- [29] B. G. Fabian Caba Heilbron, Victor Escorcia and J. C. Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *CVPR*, pages 961–970, 2015.
- [30] P. Das, C. Xu, R. F. Doell, and J. J. Corso. A thousand frames in just a few words: Lingual description of videos through latent topics and sparse object stitching. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2634–2641, 2013.
- [31] J. Romero, D. Tzionas, and M. J. Black. Embodied hands: Modeling and capturing hands and bodies together. *ACM Transactions on Graphics, (Proc. SIGGRAPH Asia)*, 36(6), Nov. 2017.
- [32] M. Loper, N. Mahmood, J. Romero, G. Pons-Moll, and M. J. Black. Smpl: A skinned multi-person linear model. *ACM transactions on graphics (TOG)*, 34(6):1–16, 2015.
- [33] G. Pavlakos, V. Choutas, N. Ghorbani, T. Bolkart, A. A. A. Osman, D. Tzionas, and M. J. Black. Expressive body capture: 3D hands, face, and body from a single image. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 10975–10985, 2019.
- [34] J. Wang, F. Mueller, F. Bernard, S. Sorli, O. Sotnychenko, N. Qian, M. A. Otaduy, D. Casas, and C. Theobalt. Rgb2hands: real-time tracking of 3d hand interactions from monocular rgb video. *ACM Transactions on Graphics (TOG)*, 39(6):1–16, 2020.
- [35] A. Kanazawa, M. J. Black, D. W. Jacobs, and J. Malik. End-to-end recovery of human shape and pose. *CoRR*, abs/1712.06584, 2017. URL <http://arxiv.org/abs/1712.06584>.

- [36] Y. Rong, T. Shiratori, and H. Joo. Frankmocap: A monocular 3d whole-body pose estimation system via regression and integration. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops*, pages 1749–1759, October 2021.
- [37] S. Han, B. Liu, R. Wang, Y. Ye, C. D. Twigg, and K. Kin. Online optical marker-based hand tracking with deep labels. *ACM Transactions on Graphics (TOG)*, 37(4):1–10, 2018.
- [38] V. Kumar and E. Todorov. Mujoco haptix: A virtual reality system for hand manipulation. In *2015 IEEE-RAS 15th International Conference on Humanoid Robots (Humanoids)*, pages 657–663, 2015. doi:10.1109/HUMANOIDS.2015.7363441.
- [39] S. Li, X. Ma, H. Liang, M. Görner, P. Ruppel, B. Fang, F. Sun, and J. Zhang. Vision-based teleoperation of shadow dexterous hand using end-to-end deep neural network. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 416–422. IEEE, 2019.
- [40] A. Handa, K. Van Wyk, W. Yang, J. Liang, Y.-W. Chao, Q. Wan, S. Birchfield, N. Ratliff, and D. Fox. Dexpivot: Vision-based teleoperation of dexterous robotic hand-arm system. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 9164–9170, 2020. doi:10.1109/ICRA40945.2020.9197124.
- [41] L. Shao, T. Migimatsu, Q. Zhang, K. Yang, and J. Bohg. Concept2robot: Learning manipulation concepts from instructions and human demonstrations. *The International Journal of Robotics Research*, 40(12-14), 2021.
- [42] A. S. Chen, S. Nair, and C. Finn. Learning generalizable robotic reward functions from” in-the-wild” human videos. *arXiv preprint arXiv:2103.16817*, 2021.
- [43] S. Bahl, A. Gupta, and D. Pathak. Human-to-robot imitation in the wild. *RSS*, 2022.
- [44] K. Schmeckpeper, O. Rybkin, K. Daniilidis, S. Levine, and C. Finn. Reinforcement learning with videos: Combining offline observations with interaction. *arXiv preprint arXiv:2011.06507*, 2020.
- [45] P. Sharma, D. Pathak, and A. Gupta. Third-person visual imitation learning via decoupled hierarchical controller. *arXiv preprint arXiv:1911.09676*, 2019.
- [46] L. Smith, N. Dhawan, M. Zhang, P. Abbeel, and S. Levine. Avid: Learning multi-stage tasks via pixel-level translation of human videos. In *RSS*, 2020.
- [47] S. Young, D. Gandhi, S. Tulsiani, A. Gupta, P. Abbeel, and L. Pinto. Visual imitation made easy. *arXiv preprint arXiv:2008.04899*, 2020.
- [48] J. Lee and M. S. Ryoo. Learning robot activities from first-person human videos using convolutional future regression. In *CVPR Workshops*, pages 1–2, 2017.
- [49] N. Das, S. Bechtel, T. Davchev, D. Jayaraman, A. Rai, and F. Meier. Model-based inverse reinforcement learning from visual demonstrations. *arXiv preprint arXiv:2010.09034*, 2020.
- [50] J. Pari, N. Muhammad, S. P. Arunachalam, L. Pinto, et al. The surprising effectiveness of representation learning for visual imitation. *arXiv preprint arXiv:2112.01511*, 2021.
- [51] S. Dasari, J. Wang, J. Hong, S. Bahl, Y. Lin, A. S. Wang, A. Thankaraj, K. S. Chahal, B. Calli, S. Gupta, et al. Rb2: Robotic manipulation benchmarking with a twist. In *NeurIPS Datasets and Benchmarks Track (Round 2)*, 2021.
- [52] A. J. Ijspeert, J. Nakanishi, H. Hoffmann, P. Pastor, and S. Schaal. Dynamical movement primitives: Learning attractor models for motor behaviors. *Neural Computation*, 2013.
- [53] M. Prada, A. Remazeilles, A. Koene, and S. Endo. Dynamic movement primitives for human-robot interaction: Comparison with human behavioral observation. In *International Conference on Intelligent Robots and Systems*, 2013.
- [54] S. Schaal. Dynamic movement primitives—a framework for motor control in humans and humanoid robotics. In *Adaptive motion of animals and machines*. Springer, 2006.

- [55] P. Pastor, H. Hoffmann, T. Asfour, and S. Schaal. Learning and generalization of motor skills by learning from demonstration. In *ICRA*, 2009.
- [56] A. Sivakumar, K. Shaw, and D. Pathak. Robotic telekinesis: Learning a robotic hand imitator by watching humans on youtube, 2022.
- [57] Allegro hand. <https://www.wonikrobotics.com/research-robot-hand>.
- [58] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *CoRR*, abs/1512.03385, 2015. URL <http://arxiv.org/abs/1512.03385>.
- [59] xarm6 by ufactory. <https://www.ufactory.cc/xarm-collaborative-robot>.
- [60] M. Contributors. Openmmlab’s next generation video understanding toolbox and benchmark. <https://github.com/open-mmlab/mmaaction2>, 2020.
- [61] Z. Cao, G. Hidalgo, T. Simon, S.-E. Wei, and Y. Sheikh. Openpose: realtime multi-person 2d pose estimation using part affinity fields. *IEEE transactions on pattern analysis and machine intelligence*, 43(1):172–186, 2019.
- [62] M. A. Fischler and R. C. Bolles. Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography. *Commun. ACM*, 24(6):381–395, jun 1981. ISSN 0001-0782. doi:10.1145/358669.358692. URL <https://doi.org/10.1145/358669.358692>.
- [63] J. L. Schönberger, E. Zheng, M. Pollefeys, and J.-M. Frahm. Pixelwise View Selection for Unstructured Multi-View Stereo. In *European Conference on Computer Vision (ECCV)*, 2016.
- [64] C. Campos, R. Elvira, J. J. G. Rodríguez, J. M. Montiel, and J. D. Tardós. Orb-slam3: An accurate open-source library for visual, visual–inertial, and multimap slam. *IEEE Transactions on Robotics*, 37(6):1874–1890, 2021.
- [65] X. Zhou, R. Girdhar, A. Joulin, P. Krähenbühl, and I. Misra. Detecting twenty-thousand classes using image-level supervision. In *ECCV*, 2022.
- [66] S. F. Bhat, I. Alhashim, and P. Wonka. Adabins: Depth estimation using adaptive bins. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4009–4018, 2021.
- [67] A. Kumar, A. Zhou, G. Tucker, and S. Levine. Conservative q-learning for offline reinforcement learning. *Advances in Neural Information Processing Systems*, 33:1179–1191, 2020.
- [68] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [69] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16000–16009, 2022.
- [70] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [71] A. Handa, K. Van Wyk, W. Yang, J. Liang, Y.-W. Chao, Q. Wan, S. Birchfield, N. Ratliff, and D. Fox. Dexpilot: Vision-based teleoperation of dexterous robotic hand-arm system. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 9164–9170. IEEE, 2020.
- [72] M. I. Takuma Seno. d3rlpy: An offline deep reinforcement library. In *NeurIPS 2021 Offline Reinforcement Learning Workshop*, December 2021.

## A Videos

Task videos, performance videos, data collection and example of internet videos can be found at: <https://videodex.github.io>

## B Additional Ablations

**Comparing Effects of Actions, Visual and Physical Priors:** Firstly, we ran an ablation where we pertained a policy on human videos performing the place task and finetune it on the uncover task (using robot data). Similarly, we pretrained a policy on Uncover and finetuned on place. The results are in the below table under VideoDex-Transfer. We see that for both tasks the performance degrades slightly, especially in the place task. We also train by adding noise to the demonstration trajectories, by adding two different levels of Gaussian noise with standard deviation being 0.01 and 0.05, shown as VideoDex-Noise-0.01 and VideoDex-Noise-0.05. We find that adding more noise definitely hurts the performance of the method. We also train ResNet18 [58] features initialized from ImageNet [70] training instead of the R3M [6] features, and the results in VideoDex-ImageNet. We can see that performance drops off, which indicates that the visual priors are important. Note that all of the reported numbers are on test objects. We present the results in Table 6.

## C Retargeting Details

We first retarget human videos from Epic-Kitchens [2]. Specifically, we use the new data (refresher) from their GoPro Hero 7 Black. We retarget video clips of humans completing tasks that are similar to the robot task. These clips are on average 5-10 seconds each, depending on the task.

**Wrist in Camera frame** The goal of Perspective-n-Point is to estimate the pose of the calibrated camera given a set of N 3D points in the world and their corresponding 2D point projections in the image. First the camera must be calibrated. To do this, we use COLMAP [63] on a set of videos. It tracks keypoints through frames and estimates the calibration from the internet videos. We find these camera intrinsics for the GoPro:

$$\begin{bmatrix} 2304.002572862 & 0 & 960 \\ 0 & 2304.002572862 & 540 \\ 0 & 0 & 1 \end{bmatrix}$$

Using this calibration, we can now complete the Perspective-n-Point process. We are given two sets of points, 16 points in 3D on the hand model in the model’s frame  $[X_w, Y_w, Z_w, 1]^t$ , and another

Transforms	Description	Method
$M_{C_t}^{Wrist}$	Wrist in each Camera	FrankMocap + PnP
$M_{C_1}^{C_t}$	Track Moving Camera	IMU/ORBSLAM
$M_{World}^{C_1}$	Make Camera parallel to Ground	IMU/Stabilization Sensor
$T_{Robot}^{World}$	Rescale and Reorient for Robot	Heuristic
$M_{Robot}^{Wrist}$	$T_{Robot}^{World} \cdot M_{World}^{C_1} \cdot M_{C_1}^{C_t} \cdot M_{C_t}^{Wrist}$	

Table 5: Transformations required to calculate wrist in robot frame from passive videos to use in learning. M denotes a transformation matrix, where T is a general transformation

Method/Task	Place	Uncover
VideoDex-Noise-0.01	0.55	0.87
VideoDex-Noise-0.05	0.50	0.60
VideoDex-ImageNet	0.40	0.62
VideoDex-Transfer	0.60	0.87
VideoDex-Original	0.70	0.90

Table 6: Ablations that compare effects of different action, visual and physical priors, as well as seeing how pretraining on different data transfers to other tasks.

set of 16 2D points in image frame  $[u, v, 1]^t$ :

$$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \begin{bmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} r_{11} & r_{12} & r_{13} & t_x \\ r_{21} & r_{22} & r_{23} & t_y \\ r_{31} & r_{32} & r_{33} & t_z \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix}$$

We use the OpenCV3 solvePnPRANSAC to complete this calculation. This implementation ensures that the process is resilient to erroneous detections.

**Camera in First Camera frame** In the SLAM section, the goal is to track the camera through the video. This is required to compensate for the movement of the camera. We use this on a selection of videos where we found the camera to move significantly.

We start the SLAM process two seconds before the action clip begins. We mark the start of the action’s frame as the first frame, run it through SLAM and then recover the trajectory of the camera through the entire clip. Specifically we recover the transformation:  $M_{C_1}^{C_t}$

The process of monocular SLAM is only valid up to a scale factor. Although Epic Kitchens has noisy accelerometer and gyro information from the camera’s sensors, we do not use this data to disambiguate this scale factor throughout the duration of the video clip.

ORB\_SLAM3 [64] only evaluates real-time video going forward through time and does not recalculate prior poses from future information. While this seems imperfect for this purpose, we find that the results are satisfactory for our purpose. In our setup, we are using these retargeted videos as a prior for learning. These retargeted trajectories are not used directly on the robot so they do not need complete accuracy. The more important characteristic is speed. COLMAP [63] can take hours to process larger video clips, but ORB\_SLAM3 [64] can complete this process faster than real time. This enables us to use many videos as an action prior for the robot behavior.

**Camera Parallel to Ground** Now that we have the trajectory in the  $C_1$  frame after SLAM and PnP, we still are missing some key transformations to get into the robot frame. First, the  $C_1$  is not always upright compared to gravity, but the robot always is. If we have a vector normal to the ground either from the synthesized pseudo-depth map from the original Videodex method or privileged information from an accelerometer (accelerometer is not used in the original Videodex method) we can use:

$$\text{pitch} = \tan^{-1}(x_{Acc}/\sqrt{y_{Acc}^2 + z_{Acc}^2}) \quad (6)$$

$$\text{roll} = \tan^{-1}(y_{Acc}/\sqrt{x_{Acc}^2 + z_{Acc}^2}) \quad (7)$$

$$\text{theta} = \tan^{-1}(\sqrt{x_{Acc}^2 + y_{Acc}^2}/z_{Acc}) \quad (8)$$

The pitch and roll would be used to make the trajectory upright. The yaw is not something that is calculable this way because this rotation is around the z axis, or the direction of gravity so it isn’t detected by an accelerometer. The theta represents how far the accelerometer z axis is off from upright but is not useful to reorient the frame.

Task	Robot Demos	Objects
Pick	125	8
Rotate	140	8
Open	120	4
Cover	124	12
Uncover	145	12
Place	175	10
Push	136	14

Table 7: Left: Number of trajectories we used for each task. Robot data is collected locally using teleportation. Most of these trajectories are 5-15 seconds in length and capture the motion trajectory of the task and visual data. Right: The number of different objects we used for each task’s data collection. In our testing, we show generalization outside of this set of objects.

**Accelerometer Robot Reorientation** There’s book-keeping transformations that must be included to rotate everything into the same frame conventions. Accelerometers, like the one in the GoPro have their frame where Z is up, y is into the screen from the lens, and x is to the left if you’re looking at the screen. The camera frame has the x-axis pointing to the right from the screen point of view, the y-axis facing down, and the z-axis facing out of the lens. The robot frame has its x-axis facing out towards the table, the y-axis faces to the left from the robot point of view, and the z-axis points up. This then leads to the following results. The camera in world frame in roll, pitch, yaw using fixed axis is:  $[pitch, 0, -roll]$ . The world frame to robot frame rotation in roll, pitch, yaw using fixed axis is  $[-3.14/2, 0, -3.14/2]$  This is used to rotate the trajectories to the robot frame and is the rotation component of  $T_{Robot}^{World}$ .

**Rescaling for Robot** We must fit the trajectories from the human videos into the robot frame. The robot frame has significant workspace limits that the human does not have. Even if the human arm is smaller than the robot’s, the human can walk around whereas the robot arm cannot move from the middle of the table. We therefore center the trajectory and ensure it fits in the robot frame. This is the scaling portion of  $T_{Robot}^{World}$ .

We rescale each dimension of the arm trajectory as:

$$M_{World}^{WristN} = M_{World}^{WristN} - (\max(M_{World}^{Wrist1..N}) + \min(M_{World}^{Wrist1..N}))/2 + \text{robotWorkspaceCenter}$$

We would like to generate more similar trajectories to use in possible data augmentation. The naive method is to add gaussian noise to the trajectory. While this can be valid, it adds noise to an already noisy system. Instead we leverage the coordinate frames to create more accurate trajectories. We randomize the workspace scaling that is used by 10 percent. Additionally, we create a rotation  $M_{World}^{World}$  that rotates the initial world frame by up to 10 degrees in each fixed axis in roll pitch yaw convention. This perturbs the direction that the robot moves in its frame.

While this augmentation can be helpful with lower amounts of internet data, in our results it was not used as it led to similar results to not using data augmentation.

We interpolate the length of the trajectories using RBF basis functions. All trajectories from the internet data are rescaled to 200 datapoints. This uniformity enables efficient batch training and was used for all of the results.

**Hand Re-targeting** We use a similar approach to re-targeting as Sivakumar et al. [56] and Handa et al. [71]. Specifically, we use the the detected human hand poses using MANO Romero et al. [31] (and FrankMocap Rong et al. [36]) to match 3D keypoints between human hands and the allegro hand. Given human hand parameters  $(\beta, \theta)$ , the goal is to minimize the difference between human

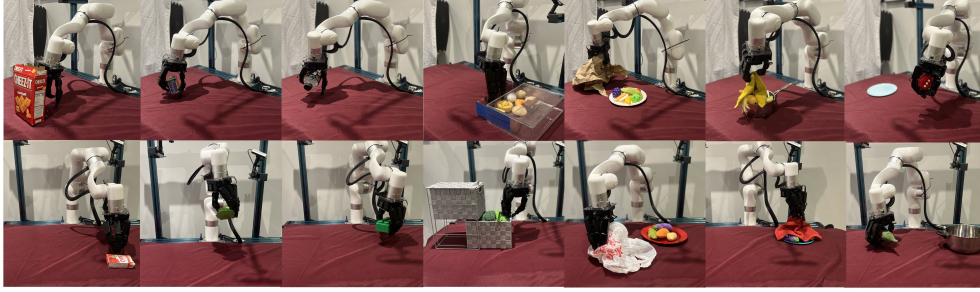


Figure 7: **Task Images.** A more detailed look at the tasks completed by VideoDex: push, pick, rotate, open, cover, uncover and place. and our website at <https://videodex.github.io> for further details.

	Pick		Rotate		Open		Cover		Uncover		Place		Push	
	train	test	train	test	train	test	train	test	train	test	train	test	train	test
BC-NDP [14]	0.64 ± 0.11	0.38 ± 0.13	0.94 ± 0.06	0.56 ± 0.13	0.90 ± 0.10	0.60 ± 0.16	0.78 ± 0.15	0.58 ± 0.15	0.88 ± 0.13	0.82 ± 0.12	0.70 ± 0.15	0.35 ± 0.11	1.00 ± 0.00	0.71 ± 0.13
BC-Open[51]	0.50 ± 0.12	0.44 ± 0.13	0.72 ± 0.11	0.38 ± 0.13	0.80 ± 0.13	0.40 ± 0.16	0.44 ± 0.18	0.58 ± 0.15	1.00 ± 0.00	0.91 ± 0.09	0.40 ± 0.16	0.25 ± 0.10	1.00 ± 0.00	0.93 ± 0.07
BC-RNN [51]	0.56 ± 0.12	0.31 ± 0.12	0.78 ± 0.10	0.50 ± 0.13	0.90 ± 0.10	0.50 ± 0.17	0.56 ± 0.18	0.42 ± 0.15	0.88 ± 0.13	0.75 ± 0.13	0.70 ± 0.15	0.50 ± 0.11	1.00 ± 0.00	1.00 ± 0.00
VideoDex	0.81 ± 0.09	0.75 ± 0.11	0.89 ± 0.08	0.69 ± 0.12	0.90 ± 0.10	0.80 ± 0.13	0.78 ± 0.15	0.67 ± 0.14	1.00 ± 0.00	0.90 ± 0.10	0.90 ± 0.10	0.70 ± 0.11	1.00 ± 0.00	1.00 ± 0.00

Table 8: We present the variance of train objects and test objects for Videodex and baselines described above. See the main paper for the mean results.

and robot keypoints: Human  $v_i^h$  and robot  $v_i^r$ . The robot keypoints are a function of robot joint pose:  $q$ . This is done by the implicit energy function ( $c_i$  are scale hyperparameters):

$$E_\pi((\beta_h, \theta_h), q) = \sum_i \|v_i^h - (c_i \cdot v_i^r)\|_2^2 \quad (9)$$

This is inefficient to compute in real time, thus similarly to Sivakumar et al. [56], we distill this into a single neural network

$$f_{\text{hand}}((\beta_h, \theta_h)) = \hat{q}$$

This network learns to minimize the energy function  $E_\pi$  described above and is trained by observing internet videos. The hand retargeting setup can be seen in Figures 3 and 4 of the main paper.

## Task

The tasks that were completed are Pick, Rotate, Open, Cover, Uncover, Place, Push. In pick, the task is to pickup objects off the table, or a plate/pan. Rotate involves turning an object in place. Open involves opening a drawer. Cover and uncover involve putting on or removing some form of cloth (dish, rubber, paper or plastic) on or from a plate. For place, the robot has to pickup an object and drop it in a plate or pot/pan. For push, the robot has to poke the object with its fingers. These tasks can be seen in Figure 7. We used videos from Epic Kitchens [2] that were as close as possible to these tasks, and doing similar types of actions. More details can be found in Table 9.

## D Learning Pipeline Details

**Learning Setup** For our approach, we use the ResNet18 from R3M [6] weights as the visual backbone. This produces an intermediate feature vector of size 512. This is processed with a 2 layer MLP with a hidden dimension of 512. The visual features are concatenated with the starting hand and wrist pose. We employ two such MLPs, one for the hand and wrist trajectories. These are then processed with an NDP [13]. The NDP processes the input with a single hidden layer to project it into the desired size (parameters  $W$  and  $g$ ). For more information we point the readers to Bahl et al. [13]. We use the implementation from Dasari et al. [51]. We use standard data augmentations from Pytorch. Specifically, we use RandomResizeCrop from a scale of 0.8 to 1.0. We use RandomGrayscale with a probability of 0.05. We use ColorJitter with a brightness of 0.4, contrast of 0.3, saturation of 0.3 and hue of 0.3. Finally, we normalize the RGB values around the typical mean and standard deviation for color images:  $\mu = (0.485, 0.456, 0.406)$   $\sigma = (0.229, 0.224, 0.225)$ . For different baselines, we used the same backbone (R3M [6]) as our method. We use the same architecture style as well, with the visual features being processed by both a wrist and hand stream. Finally, all network sizes are the same or very similar. We describe our hyperparameters in Table 9.



Parameter	Value
Learning Rate	$1 \times 10^{-3}$
Batch Size	32
Training Demonstrations Per Task	120-175
Human Videos Per Task	350 (Cover/Unicver, Rotate, Push) - 2500 (Open, Pick, Place)
Trajectory Length Human Videos	200 (rescaled)
NDP [13] Basis Functions $N$	300
NDP [13] Global Parameter $\alpha$	15

Table 9: Parameter List

## E Experimental Setup

We collect data using a dexterous hand robotic teleoperation setup [56, 71]. A trained operator stands in front of the camera within view of the robot and operates the system in real-time to collect demonstrations with a trained, uniform style. Another manager stands by. The goal of this manager is to place items on the table for manipulation, randomize locations and types of objects, to start and stop demonstrations for the operator and manage the robot system. We collect about 120-175 demonstrations per task. See table 5 for details.

## F Hardware Details

Our hardware setup consists of an LEAP 16 DOF Hand and an XArm 6 manipulator (from Ufactory). The hand is mounted on the wrist of the XArm. To collect data we use a similar teleoperation system as provided by Sivakumar et al. [56] and Handa et al. [71]. We use Intel Realsense D415 cameras to collect human teleoperation and robot videos. We use four NVIDIA RTX 2080TI's for training the policy and running the teleoperated system. We experiment with the Allegro Hand and use it to collect some teleoperated demonstration data, but we find it to be very unreliable and break many times. The motors also quickly overheat and are weak in practice. Therefore, to alleviate these issues we use the LEAP Hand for collecting the final results.

### External Codebases

We use the following different external codebases for our pipeline:

- Human body and hand detection: FrankMocap [36], <https://github.com/facebookresearch/frankmocap>
- NDP and Behavior Cloning [51] code from [https://github.com/AGI-Labs/robot\\_baselines](https://github.com/AGI-Labs/robot_baselines)
- Code for R3M [6] <https://github.com/facebookresearch/r3m>
- CQL baseline from Takuma Seno [72] (<https://github.com/takuseno/d3rlpy>)
- COLMAP from Schönberger et al. [63] (<https://github.com/colmap/colmap>)
- ORBSLAM3 from Campos et al. [64] ([https://github.com/UZ-SLAMLab/ORB\\_SLAM3](https://github.com/UZ-SLAMLab/ORB_SLAM3))
- GoPro Metadata Extractor from (<https://github.com/JuanIrache/gpmf-extract>)
- Rigid Transform class from ([https://github.com/BerkeleyAutomation/autolab\\_core/blob/master/autolab\\_core/rigid\\_transformations.py](https://github.com/BerkeleyAutomation/autolab_core/blob/master/autolab_core/rigid_transformations.py))
- PnP from (<https://opencv.org/>)