# One-Shot Hierarchical Imitation Learning of Compound Visuomotor Tasks

Tianhe Yu[1], Pieter Abbeel[1], Sergey Levine[1], Chelsea Finn[1]

*Abstract*— **We consider the problem of learning multi-stage vision-based tasks on a real robot from a single video of a human performing the task, while leveraging demonstration data of subtasks with other objects. This problem presents a number of major challenges. Video demonstrations without teleoperation are easy for humans to provide, but do not provide any direct supervision. Learning policies from raw pixels enables full generality but calls for large function approximators with many parameters to be learned. Finally, compound tasks can require impractical amounts of demonstration data, when treated as a monolithic skill. To address these challenges, we propose a method that learns both how to learn primitive behaviors from video demonstrations and how to dynamically compose these behaviors to perform multi-stage tasks by "watching" a human demonstrator. Our results on a simulated Sawyer robot and real PR2 robot illustrate our method for learning a variety of order fulfillment and kitchen serving tasks with novel objects and raw pixel inputs. Video results are linked at `https://sites.google.com/view/one-shot-hil`.**

## I. INTRODUCTION

Humans have a remarkable ability to imitate complex behaviors from just *watching* another person. In contrast, robots require a human to physically guide or teleoperate the robot's body [1]–[5] in order to learn from demonstrations. Furthermore, people can effectively learn behaviors from just a single demonstration, while robots often require substantially more data [5], [6]. While prior work has made progress on imitating manipulation primitives using raw video demonstrations from humans [7]–[10], handling more complex, compound tasks presents an additional challenge. When skill sequences become temporally extended, it becomes impractical to treat a sequence of skills as a single, monolithic task, as the full sequence has a much longer time horizon than individual primitives and learning requires significantly more data. In this paper, we consider the following question: can we leverage the compositional structure underlying compound tasks to effectively learn temporally-extended tasks from a single video of a human demonstrator?

A number of prior works aim to learn temporally extended skills using demonstration data that is labeled based on the particular primitive executed or using pre-programmed primitives [11]–[14]. However, when the human demonstrations are provided as raw videos, and robot must also handle raw visual observations, it is difficult to employ conventionally methods such as low-dimensional policy representations [15], [16] or changepoint detection [17]. In this paper, we consider a problem setting of learning to perform multi-stage tasks through imitation where the robot must map
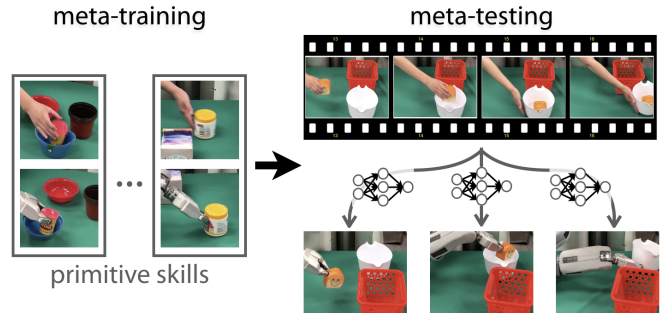
[1] Berkeley AI Research, UC Berkeley Computer Science. Correspondence to `tianheyu@cs.stanford.edu`

Fig. 1. A robot learns and composes convolutional neural network policies for performing a multi-stage manipulation task by "watching" a human perform the compound task (right), by leveraging data from primitive skills with other objects (left).

raw image observations to actions, and the demonstration is provided via an unsegmented raw video of a person performing the entire task. To approach this problem, the key idea in this work is to leverage meta-learning, where the robot uses *previous* data of primitive skills to learn how to imitate humans performing multi-stage skills. In particular, the meta-training set of previous *primitive* skills consists of both videos of humans and teleoperated episodes, while the new *multi-stage* tasks seen at meta-test time are only specified with a single video of a human, without teleoperation (see Figure 1). Hence, our goal is to both learn primitive behaviors and to compose them automatically from a single video of a human performing the new compound task.

Our approach adapts its policy online over the course of the multi-stage task as it "watches" the video of the human. To accomplish this, we build models that can recognize the progress towards completion of the current primitive (or the primitive's 'phase'), and integrate these models with a one-shot imitator to dynamically learn and compose primitives into compound tasks, from a video of a human. The phase of a primitive can be learned directly from the demonstration data of primitives, using the frame indices of each demonstration, without requiring any manual labeling. To learn a policy for an individual primitive from a video of a human performing that primitive, a building block of our approach, we use domain-adaptive meta-imitation learning [9]. All in all, as illustrated in Figure 2, our method decomposes the test-time human video into primitives using a primitive phase predictor, computes a sequence of policies for each primitive, and sequentially executes each policy until each has deemed it is complete, again utilizing a phase predictor.

The primary contribution of this work is an approach for dynamically learning and composing sequences of policies based on a single human demonstration without annotations.

Our method combines one-shot imitation of subtasks with a learned mechanism for decomposing compound task demonstrations and composing primitive skills. In our experimental results, we show that this approach can be used to dynamically learn and sequence skills from a single video demonstration provided by the user at test time. We evaluate our method on order fulfillment tasks with a simulated Sawyer arm and a kitchen serving tasks with a real PR2 robot. Our method learns to perform temporally-extended tasks from raw pixel inputs and outperforms alternative approaches.

## II. RELATED WORK

Prior methods for compound task learning from demonstrations often use demonstrations that are labeled or segmented into individual activities [12], commands [13], [14], [18], or textual descriptions [19], or assume a library of pre-specified primitives [20], [21]. Unlike many of these approaches, we use a dataset of demonstrations that are not labeled by activity or command, reducing the burden on labeling. Further, at test time, we translate directly from a human video to a sequence of policies, bypassing the intermediate grounding in language or activity labels. This "end-to-end" approach has the benefit of making it straightforward to convey motions or subtasks that are difficult to describe, such as the transitions from one subtask to another; it also does not require the human to have any knowledge of the grammar of commands. Long-horizon tasks have also been considered in one-shot imitation works [22], [23]. We consider a different setting where only demonstrations of primitives are available during meta-training, rather than demonstrations of full-length compound tasks.

A number of prior works have aimed to decompose demonstrations into primitives or skills, e.g. using change point detection [17], latent variable models [24]–[26], mixture of experts [27], or transition state clustering [28]. Like these prior works, we build upon the notion that each primitive has its own policy, aiming to construct a sequence of policies, each with a learned termination condition. Similar to Manschitz et al. [29], we learn the termination condition by predicting the phase of a primitive skill. Unlike these approaches, we consider the problem of learning an end-to-end visuomotor policy (pixels to end-effector velocities) from a single video of a human performing a task, while leveraging visual demonstration data from primitives performed with previous objects.

## III. PRELIMINARIES

To learn how to learn and compose primitive skills into multi-stage tasks, we need to first learn how to learn a primitive skill by imitating a human demonstration. To do so, our approach uses prior work on domain-adaptive meta-learning (DAML) [9] that learns how to infer a policy from a single human demonstration. DAML is an extension of the model-agnostic meta-learning algorithm (MAML) [30]. In this section, we present the overview meta-learning problem, discuss both MAML and DAML, and introduce notation.

A primary goal of many meta-learning algorithms is to learn new tasks with a small amount of data. To achieve this, these methods learn to efficiently learn many meta-training tasks, such that, when faced with a new meta-test task, it can be learned efficiently. Note that meta-learning algorithms assume that the meta-training and meta-test tasks are sampled from the same distribution $p(\mathcal{T})$, and that there exists a common structure among tasks, such that learning this structure can lead to fast learning of new tasks (referred as few-shot generalization). Hence, meta-learning corresponds to structure learning.

MAML aims to learn the shared structure across tasks by learning parameters of a deep network such that one or a few steps of gradient descent leads to effective generalization to new tasks. We will use $\theta$ to denote the initial model parameters and $\mathcal{L}(\theta, \mathcal{D})$ to denote the loss function of a supervised learner, where $\mathcal{D}_{\mathcal{T}}$ denotes labeled data for task $\mathcal{T}$. During meta-training, MAML samples a task $\mathcal{T}$ and datapoints from $\mathcal{D}_{\mathcal{T}}$, which are randomly partitioned into two sets, $\mathcal{D}_{\mathcal{T}}^{\mathrm{tr}}$ and $\mathcal{D}_{\mathcal{T}}^{\mathrm{val}}$. We will assume that there are $K$ examples in $\mathcal{D}_{\mathcal{T}}^{\mathrm{tr}}$. MAML optimizes for model parameters $\theta$ such that one or a few gradient steps on $\mathcal{D}_{\mathcal{T}}^{\mathrm{tr}}$ results in good performance on $\mathcal{D}_{\mathcal{T}}^{\mathrm{val}}$. Concretely, the MAML objective is:

$$\min_{\theta} \sum_{\mathcal{T} \sim p(\mathcal{T})} \mathcal{L}(\theta - \alpha \nabla_\theta \mathcal{L}(\theta, \mathcal{D}_{\mathcal{T}}^{\mathrm{tr}}), \mathcal{D}_{\mathcal{T}}^{\mathrm{val}})$$
$$= \min_{\theta} \sum_{\mathcal{T} \sim p(\mathcal{T})} \mathcal{L}(\phi_{\mathcal{T}}, \mathcal{D}_{\mathcal{T}}^{\mathrm{val}}).$$

where $\phi_{\mathcal{T}}$ corresponds to the updated parameters and $\alpha$ is a step size of the gradient update. Moving forward, we will refer to the inner loss function as the *adaptation objective* and the outer objective as the *meta-objective*. At meta-test time, in order to infer the updated parameters for a new, held-out task $\mathcal{T}_{\mathrm{test}}$, MAML runs gradient descent with respect to $\theta$ using $K$ examples drawn from $\mathcal{T}_{\mathrm{test}}$:

$$\phi_{\mathcal{T}_{\mathrm{test}}} = \theta - \alpha \nabla_\theta \mathcal{L}(\theta, \mathcal{D}_{\mathcal{T}_{\mathrm{test}}}^{\mathrm{tr}}).$$

DAML applied the MAML algorithm to the domain-adaptive one-shot imitation learning setting; DAML aims to learn how to learn from a video of a human, using teleoperated demonstrations for evaluating the meta-objective. Essentially, DAML learns to translate from a video of a human performing a task to a policy that performs that task. Unlike the standard supervised meta-learning setting, a video of a human is a form of weak supervision – it contains enough information to communicate the task, but does not contain direct label supervision. To allow a robot to learn from a video of a human and handle the domain shifts between the human and the robot, DAML additionally learns the *adaptation objective* denoted as $\mathcal{L}_\psi$ along with the initial parameters $\theta$. The meta-objective is a mean-squared error behavioral cloning loss denoted as $\mathcal{L}_{\mathrm{BC}}$. $\mathcal{L}_\psi$ can be viewed as a learned critic such that running gradient descent on $\mathcal{L}_\psi$ can produce effective adaptation. Specifically, the DAML meta-objective can be formulated as follows:

$$\min_{\theta, \psi} \sum_{\mathcal{T} \sim p(\mathcal{T})} \sum_{\mathbf{d}^h \in \mathcal{D}_{\mathcal{T}}^h} \sum_{\mathbf{d}^r \in \mathcal{D}_{\mathcal{T}}^r} \mathcal{L}_{\mathrm{BC}}(\theta - \alpha \nabla_\theta \mathcal{L}_\psi(\theta, \mathbf{d}^h), \mathbf{d}^r).$$
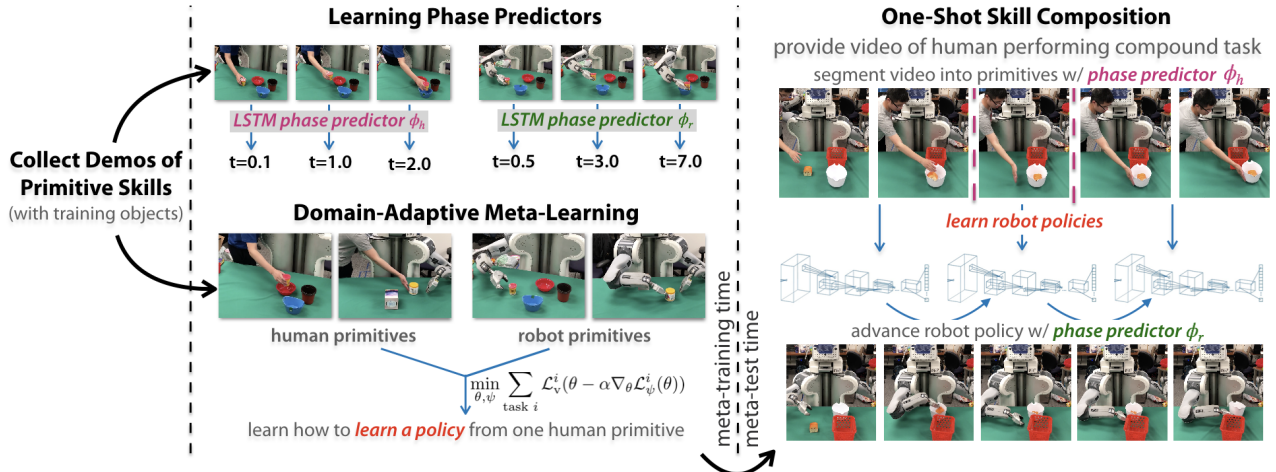
Fig. 2. After learning a phase predictor and meta-learning with human and robot demonstration primitives, the robot temporally segments a human demonstration (just a video) of individual primitives and learns to perform each segmented primitive sequentially.

where $\mathbf{d}^h$ and $\mathbf{d}^r$ are a human and robot demonstration respectively. In prior work with DAML, different tasks corresponded to different objects, hence having the robot learn to perform previous skills with novel objects. In this work, different meta-learning tasks will correspond to different primitives performed with different sets of objects. Next, we present our approach for one-shot visual skill composition.

## IV. ONE-SHOT VISUAL SKILL COMPOSITION

Our goal is for a robot to learn to perform a variety of multi-stage tasks by "watching" a video of a human perform each task. To do so, the robot needs to have some form of prior knowledge or experience, which we will obtain by using demonstration data of primitive subtasks performed with other objects. After meta-training, at meta-test time, the robot is provided with a single video of a human performing a multi-stage task, and is tasked with learning to perform the same multi-stage task in new settings where the objects have moved into different starting configurations. In this section, we define the problem assumptions and present our approach for tackling this problem setting.

### A. Problem Setting and Overview

Concretely, during an initial meta-training phase, we provide a dataset of primitive demonstrations. For each subtask with a particular set of objects (which we refer to as a primitive $\mathcal{P}_k$), we provide multiple human demonstrations, $\{\mathbf{d}_i^h\}_k$ and multiple robot demonstrations $\{\mathbf{d}_j^r\}_k$. We define a demonstration performed by a human $\mathbf{d}_k^h$ to be a sequence of images $\mathbf{o}_1^h, \ldots, \mathbf{o}_{T_i}^h$ of a human performing $\mathcal{P}_k$, and a robot demonstration $\mathbf{d}_k^r$ to be a sequence of images and actions $\mathbf{o}_1^r, \mathbf{a}_1^r \ldots, \mathbf{o}_{T_j}^r, \mathbf{a}_{T_j}^r$ of a robot performing the same primitive. Note that demonstrations have potentially different time horizons, namely it may be that $T_i \neq T_j$. The human and robot demonstration need to have a correspondence based on the objects used and primitive performed, but do not need to be aligned in any way, nor be executed with the same speed, nor have the same object positions.

After meta-training, the robot is provided with a human demonstration $\mathbf{d}^h$ of a compound task consisting of multiple primitives in sequence. The sequence of primitives seen in this video at meta-test time may involve novel objects and novel configurations thereof, though we assume that the general families of subtasks in the human demonstration (e.g., pushing, grasping, etc) are seen in the meta-training data. However, the meta-training data does not contain any composition of primitives into multi-stage tasks. After observing the provided human video, the robot is tasked with identifying primitives, learning policies for each primitive in the human's demonstration, and composing the policies together to perform the task.

In our approach, we use the meta-training data to have the robot learn how to learn an individual subtask from a video demonstration, in essence, translating a video demonstration into a policy for the subtask. Additionally, the robot needs to learn how to identify and compose learned primitives into multi-stage skills. To do so, we propose to train two models that can identify the temporal progress, or *phase*, of any human or robot primitive respectively, which can be used for segmentation of human demonstrations and for terminating a primitive being performed by the robot to move onto the next subtask. An overview of our method is in Figure 2. In the remainder of this section, we discuss how we can learn and compose primitives for completing compound tasks from single demonstrations, and how our two phase predictors can be trained.

### B. One-Shot Composition of Primitives

During meta-training, we train a human phase predictor $\phi_h$ and a robot phase predictor $\phi_r$ (as described in the following section), as well as a DAML one-shot learner $\pi_\theta$, which can learn primitive-specific robot policies $\pi_{\phi_i}$ from videos of humans performing primitives, $\mathbf{d}_i^h$. At meta-test time, the robot is provided with a video of a human completing a multi-stage task, $\mathbf{o}_1^h, \ldots, \mathbf{o}_T^h$.

We first need to decompose the multi-stage human demonstration into individual primitives. We can do so by using the human phase predictor. In particular, we feed the demonstration, frame by frame, into $\phi_h$ until $\phi_h(\mathbf{o}_{1:t}^h) > 1 - \epsilon$, indicating the end of the current primitive, $\mathbf{d}^h = \mathbf{o}_1^h, \ldots, \mathbf{o}_t^h$.

**Algorithm 1** One-Shot Composition of Primitives

---

**Require:** meta-learned initial parameters $\theta$, and adaptation objective $\mathcal{L}_\psi$

**Require:** human and robot phase predictors $\phi_h$ and $\phi_r$

**Require:** A human video for a compound task $\mathbf{o}_1^h, \dots, \mathbf{o}_T^h$

  # *Decompose human demonstration into primitives*

  Initialize predicted primitives $\mathcal{P} = \{\}$ and $t = 1, t' = 1$

  **while** $t < T$ **do**

    **if** $\phi_h(\mathbf{o}_{t':t}^h) > 1 - \epsilon$ **then**

      Append $\mathbf{d}^h = \mathbf{o}_{t':t}^h$ to $\mathcal{P}$

      $t' = t + 1$

    **end if**

    $t = t + 1$

  **end while**

  # *Compute and compose policies for each primitive*

  Initialize $t = 1$, $t' = 1$ observe $\mathbf{o}_1^r$

  **for** $\mathbf{d}_i^h$ in $\mathcal{P}$ **do**

    Compute policy parameters $\phi_i = \theta - \alpha \nabla_\theta \mathcal{L}_\psi(\theta, \mathbf{d}_i^h)$

    **while** $\phi_r(\mathbf{o}_{t':t}^r) < 1 - \epsilon$ **do**

      Take one step, executing $\pi_{\phi_i}(\mathbf{o}_t^r)$ and getting $\mathbf{o}_{t+1}^r$

      $t = t + 1$

    **end while**

    $t' = t$

  **end for**

---

Then, we repeat this process starting from the following timestep $t + 1$, to iteratively determine the endpoint of each segment. At the end of this process, we are left with a sequence of primitive demonstrations, $\mathbf{d}_1^h, \mathbf{d}_2^h, \dots$, which we translate into policies $\pi_{\phi_1}, \pi_{\phi_2}, \dots$ using the one-shot imitator $\pi_\theta$. To compose these policies together, we roll-out each policy sequentially, passing the observed images $\mathbf{o}_t$ into the robot phase predictor and advancing the policy when $\phi_r(\mathbf{o}_{1:t}^r) > 1 - \epsilon$. We detail this meta-test time process in Algorithm 1.

*C. Primitive Phase Prediction*

In order to segment a video of a compound skill into a set of primitives, many prior works train on a dataset consisting of multi-stage task demonstrations and corresponding labels that indicate the primitive at each time step. We assume a dataset of demonstrations of primitives, but without any label for the particular primitive that was executed in that segment, and where the primitive demonstrations contain different objects than those seen at meta-test time. As a result, we don't need a dataset of demonstrations of temporally extended skills during training, which can reduce the burden on humans to collect many long demonstrations. However, we still need to learn both (a) how to segment human demonstrations of compound tasks, in order to learn a sequence of policies that perform the task shown at meta-test time, and (b) when to transition from one learned policy to the next. We propose a single approach that can be used to tackle both of these problems: predicting the *phase* of a primitive. In particular, given a video of a partial demonstration or execution of any primitive, we aim to predict how much temporal progress has been made towards the completion of that primitive. This form of model can be used both to segment a video of a compound task and to identify when to switch to a subsequent policy.

We train two phase predictor models, one on human demonstration data and one on robot demonstration data. Other than the data, both models have identical form and are trained in the same way. As discussed previously, the meta-training dataset is composed of human and robot demonstrations of a set of primitives. To construct supervision for training the phase predictors, we can compute the phase of particular partial video demonstration $\mathbf{o}_{1:t}$ as $\frac{t}{T}$ where $T = |\mathbf{d}|$ is the length of the full demonstration $\mathbf{d}$, which varies across demonstrations. The phase information provides label supervision that indicates how much of the primitive has been completed. Each phase predictor is trained to take as input a partial demonstration, i.e. the first $t$ frames of a demonstration, and output the phase.

Formally, we denote $\phi_h$ as the human phase predictor, which takes as input $\mathbf{o}_{i,1:t}^h = \mathbf{o}_{i,1}^h, \dots, \mathbf{o}_{i,t}^h$ and regresses $\phi_h(\mathbf{o}_{i,1:t}^h)$ to $\frac{t}{T_i}$. Similarly, $\phi_r$ is defined to be the robot phase predictor, regressing $\phi_r(\mathbf{o}_{j,1:t}^r)$ to $\frac{t}{T_j}$, where the partial demonstration $\mathbf{o}_{j,1:t}^r = \mathbf{o}_{j,1}^r, \dots, \mathbf{o}_{j,t}^r$ only uses image observations for simplicity. To handle variable length sequences as input, both models $\phi_h$ and $\phi_r$, are represented using recurrent neural networks. Both networks are trained using a mean-squared error objective $\sum_i \sum_{t=1}^{T_i} \left\| \phi(\mathbf{d}_{i,1:t}) - \frac{t}{T_i} \right\|^2$ and the Adam optimizer [31]. The phase predictors use 5 convolution layers with 64 $3 \times 3$ filters, followed by an LSTM with 50 hidden units and a linear layer to output the phase. The first convolution layer is initialized using pretrained weights from VGG-16. We use the swish non-linearity [32], layer normalization [33], and dropout [34] with probability 0.5 at each convolution layer and the LSTM.

## V. EXPERIMENTS

The goal of our experiments is to determine whether our approach can enable a robot to learn to perform compound vision-based tasks from a single video demonstration of a human performing that task, composing primitives on the fly and generalizing to new objects. Note that this is an exceptionally challenging task: in contrast to prior work on compound task learning [11]–[14], the robot must perform the task entirely end-to-end from RGB images and receives only a single unsegmented video demonstration showing a person performing the task. The meta-training data does not contain instances of the same objects that are seen in the test-time behaviors, requiring the robot to adapt its policy to each of the objects from the prior learned during meta-training. Each task requires sequencing multiple primitives, such that no single policy can perform the entire task alone. Success on this problem requires simultaneously interpreting the human demonstration to determine which objects matter for the task and how they should be manipulated, adapting the policy multiple times to the multiple stages of the task, and executing these adapted policies in the right sequence.

Our experimental set-ups involve pick-and-place primitives, push primitives, and reach primitives. Hence, our method is capable of composing these primitives in a variety of ways to form many different compound tasks. In our evaluation, we focus on two different sets of meta-test tasks:

*a) Simulated order fulfillment:* A simulated Sawyer robot must learn to pick and place a particular set of novel objects into a bin and push the bin to a specified location.

*b) PR2 kitchen serving:* In this setting, the PR2 must grasp an object, place it into the correct bowl or platter, and push one of the platters or bowls to the robot's left.

To our knowledge, no prior work proposes an approach that aims to solve the problem that we consider. Because we do not assume access to demonstrations of compound tasks during meta-training, it is not suitable to directly compare to direct one-shot imitation on compound tasks. Even if there are many compound task demonstrations for meta-training, no prior work has demonstrated one-shot imitation learning of temporally extended tasks from raw pixels. Consequently, we compare only to ablations of our method to better understand the importance of using phase prediction and using DAML over alternative options. For understanding the former, we compare to a simple alternative to using phase prediction: learning policies for every fixed-length segment of the human demonstration and advancing the policy at every timestep. This 'sliding window' approach still leverages one-shot imitation, but makes two simplifying assumptions: (a) that fixed-length windows are an appropriate representation of primitives, and (b) that the human demonstration time and robot execution time are equal. This latter assumption will be true in our simulated experiments, where the provided demonstrations are from a robot, but may be easily violated in the real world. To study the importance of using DAML, we compare to using an LSTM-based meta-learner, akin to model of Duan et al. [22] but using visual inputs. All methods use Cartesian end-effector control. The value of $\epsilon$ is set to $0.03$, which was found to work well on validation tasks. For video results, see the project website[1].

### A. Simulated Order Fulfillment

We first evaluate our method on a range of simulated order fulfillment tasks using a Sawyer robot arm in the MuJoCo physics engine [35], as illustrated in Figure 3. Across tasks, the particular objects and number of objects to be put in the bin vary. Success is defined as having the correct objects in the bin and the bin at the target. In this experiment, we consider a simplified setting, where a typical demonstration (with both images and actions) is available at meta-test time. The next section with physical robot experiments will consider learning from human videos.

There are two types of primitive demonstrations in the training dataset: picking and placing a single object into the bin, and pushing the bin to the goal. In lieu of a simulated human demonstrator, we optimize an expert policy for pick and place using proximal policy optimization (PPO) [36] and

script an expert pushing policy. The PPO expert policy uses priveleged low-level state information (i.e. the position of the target object) rather than vision, enabling us to train a single expert policy across all pick & place primitives. To gather pick & place primitives for the training dataset, we sample one target object and two or three distractors from a set of 37 types of textures where each type contains roughly 150 different textures (for a total of around 5500 textures). We also randomize the positions of the target and distractors to form 8 different demonstrations for each primitive. We use 1800 different pick-and-place primitive subtasks, along with 56 bin pushing demonstrations, for meta-learning and training the phase predictors.

We represent the DAML network with 4 convolution layers with 24 $5 \times 5$ filters, followed by 3 fully-connected layers with 200 hidden units. We use linear adaptive objectives [9] for both actions and the gripper pose, and a step size $\alpha = 0.05$. The meta-objective is the same as prior work, as is the LSTM meta-learner architecture [9].

We evaluate each approach using one-object and two-object order fulfillment tasks. For each multi-stage meta-test task, we generate one visual demonstration by temporally concatenating expert demonstrations of the two or three primitives involved in the task, as illustrated in Figure 3. We compute success averaged over 10 tasks of each type and 3 trials per task. The results, shown in Table I, indicate that both gradient-based meta-learning and phase prediction are essential for good performance. We also observe that performance drops as the task becomes longer, indicative of compounding errors that are known to be a problem when using behavioral cloning based approaches [37]. Generally, we find that the nearly all failure cases are caused by the one-shot imitation learner, primarily related to grasping, which implicitly requires precise visual object localization and precise control. Hence, future advances in one-shot visual imitation learning will likely lead to improvements in our approach as well.

|  | 1 object | 2 objects |
|---|---|---|
| sliding window (no phase prediction) | 50.0% | 16.7% |
| LSTM one-shot learner (no DAML) | 0.0% | 0.0% |
| **one-shot skill composition (ours)** | **73.3%** | **46.7%** |

TABLE I. One-shot success rate of simulated Sawyer robot performing order fulfillment from a single demonstration with comparisons

### B. PR2 Kitchen Serving

In our second experiment on a physical PR2 robot, we collect primitive demonstrations by (a) using the dataset from [9], which contains 1293, 640, 1008, and 600 robot demonstrations for placing, pushing, pick-and-place, and (b) primitives that transition from placing to pushing, plus an equal number of human demonstrations for the above four primitives. We focus on two particular tasks for this experiment: picking an object, placing it into a particular container, and then pushing either the container with the object or a different container toward a goal position. The success metric for this task is that the object is placed into the container and the container is pushed toward the robot's left gripper.
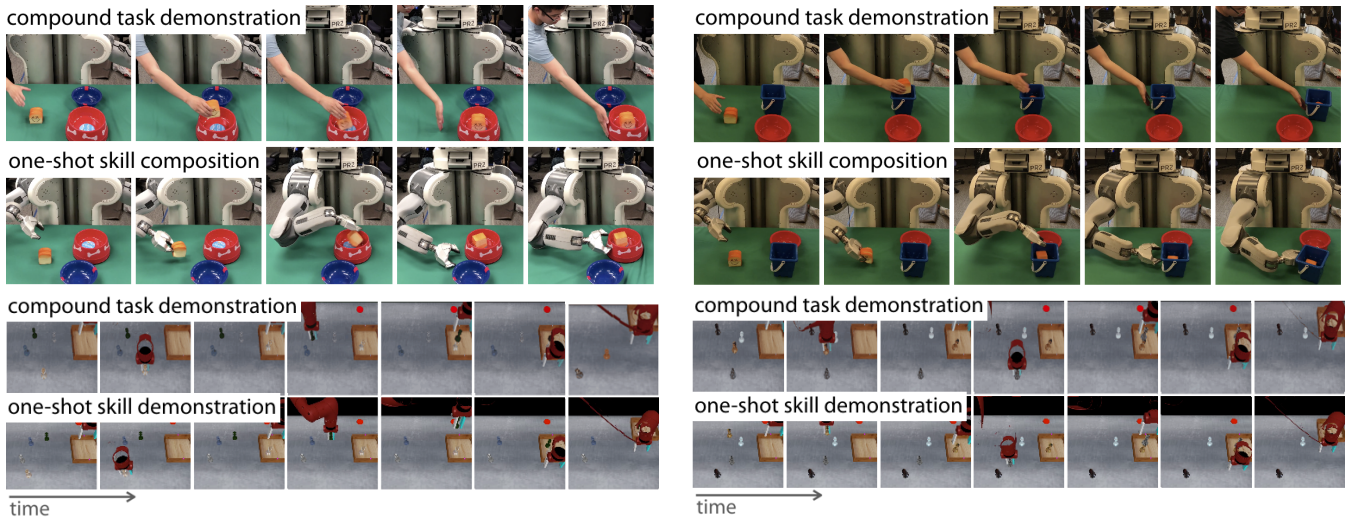
Fig. 3. Qualitative results of the experiments on the physical PR2 (top) and simulated Sawyer (bottom). The top right shows an example failure case, where the robot successfully picks and places the object into the correct container, but incorrectly pushes between the two objects. The other examples all illustrate successful learning and skill composition of the demonstrated compound task from a single demonstration.
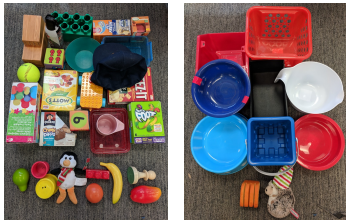


Fig. 4. A subset of the meta-training objects (left) and meta-test objects (right) used in the PR2 experiments. The top of each picture shows placing and pushing target objects and the bottom shows picked objects.

|  | same target | different target |
|---|---|---|
| sliding window (no phase prediction) | 0/20 | 0/20 |
| **one-shot skill composition (ours)** | **10/20** | **7/20** |

TABLE II. Success rate of PR2 robot performing pick & place then push task from a video of a human. Evaluated using novel, unseen objects.

## VI. DISCUSSION

We presented an approach for one-shot learning and com-position policies for achieving compound, multi-stage tasks from raw pixel inputs based on a single video of a human performing the task. Our approach leverages demonstrations from previous primitive skills in order to learn to identify the end of a primitive and to meta-learn policies for primitives. At meta-test time, our approach learns multi-stage tasks by decomposing a human demonstration into primitives, learning policies for each primitive, and composing the policies online to execute the full compound task. Our experiments indicated that our approach can successfully learn and compose convolutional neural network policies for order fulfillment tasks on a simulated Sawyer arm and a real PR2 arm, all from raw pixel input.

In future work, we hope to improve upon the performance of our approach. To do so, it will be important to improve the performance of one-shot imitation learning methods (a subcomponent of our approach) and potentially incorporate reinforcement learning or other forms of online feedback, such as DAgger [37], to overcome compounding errors. Another interesting direction for future work is to consider more unstructured demonstration data, such as human and robot demonstrations of temporally-extended tasks that are from distinct sources. Unlike many prior works that have considered the related problem of unsupervised segmentation of compound task demonstrations into primitives, this prob-lem is significantly more challenging, since it also involves learning a rough alignment between skills performed by humans and those performed by robots. But if possible, the method would provide the ability to scale to large, unlabeled datasets and may remove the need for an expert to decide the set of skills that constitute primitives during training.

We follow the same DAML architectures as in prior work [9] with two small differences. Instead of fully-connected layers, the model uses temporal convolution lay-ers; and rather than using a mixture density network at the output, the architecture uses a discretized action space and a cross-entropy loss with each action dimension independently discretized over 50 bins.

For evaluation, we use 10 novel bowls and containers and 5 novel items to be 'served' (see Figure 4). For each task, we collect a single human demonstration and evaluate two trials of the robot's task execution, as illustrated in Figure 3. We summarize the results in Table II. Since this is a challenging task as the robot needs to complete a sequence of control primitives and any small drift during imitation could lead to failure, our approach can only succeed 10 out of 20 trials when pick-and-placing the object into the target the container and pushing the correspond container, and 7 out of 20 trials when pushing the other container without the placed object. However, the naïve sliding window approach never successfully completes the entire sequence of primitives, indicating that phase prediction is an important aspect of our approach. As before, we observe that most of the failure cases of our approach hinge on the one-shot imitation learning failing to perform the primitive in the segmented video (see example in Figure 3), suggesting that future advances in one-shot imitation from human videos will lead to improved performance.

REFERENCES

[1] P. Pastor, L. Righetti, M. Kalakrishnan, and S. Schaal, "Online movement adaptation based on previous sensor experiences," in *International Conference on Intelligent Robots and Systems (IROS)*, 2011.

[2] B. Akgun, M. Cakmak, J. W. Yoo, and A. L. Thomaz, "Trajectories and keyframes for kinesthetic teaching: A human-robot interaction perspective," in *International Conference on Human-Robot Interaction*, 2012.

[3] S. Calinon, P. Evrard, E. Gribovskaya, A. Billard, and A. Kheddar, "Learning collaborative manipulation tasks by demonstration using a haptic interface," in *International Conference on Advanced Robotics (ICAR)*, 2009.

[4] R. Rahmatizadeh, P. Abolghasemi, L. Bölöni, and S. Levine, "Vision-based multi-task manipulation for inexpensive robots using end-to-end learning from demonstration," *arXiv:1707.02920*, 2017.

[5] T. Zhang, Z. McCarthy, O. Jow, D. Lee, K. Goldberg, and P. Abbeel, "Deep imitation learning for complex manipulation tasks from virtual reality teleoperation," *International Conference on Robotics and Automation (ICRA)*, 2018.

[6] R. Rahmatizadeh, P. Abolghasemi, A. Behal, and L. Bölöni, "Learning real manipulation tasks from virtual demonstrations using lstm," *AAAI Conference on Artificial Intelligence (AAAI)*, 2018.

[7] Y. Liu, A. Gupta, P. Abbeel, and S. Levine, "Imitation from observation: Learning to imitate behaviors from raw video via context translation," *International Conference on Robotics and Automation (ICRA)*, 2018.

[8] P. Sermanet, C. Lynch, J. Hsu, and S. Levine, "Time-contrastive networks: Self-supervised learning from multi-view observation," *arXiv:1704.06888*, 2017.

[9] T. Yu, C. Finn, A. Xie, S. Dasari, T. Zhang, P. Abbeel, and S. Levine, "One-shot imitation from observing humans via domain-adaptive meta-learning," *Robotics: Science and Systems (R:SS)*, 2018.

[10] B. Stadie, P. Abbeel, and I. Sutskever, "Third-person imitation learning," *International Conference on Learning Representations (ICLR)*, 2017.

[11] S. Manschitz, J. Kober, M. Gienger, and J. Peters, "Learning movement primitive attractor goals and sequential skills from kinesthetic demonstrations," *Robotics and Autonomous Systems*, vol. 74, pp. 97–107, 2015.

[12] Y. Yang, Y. Li, C. Fermüller, and Y. Aloimonos, "Robot learning manipulation action plans by" watching" unconstrained videos from the world wide web." in *AAAI Conference on Artificial Intelligence (AAAI)*, 2015.

[13] D. Xu, S. Nair, Y. Zhu, J. Gao, A. Garg, L. Fei-Fei, and S. Savarese, "Neural task programming: Learning to generalize across hierarchical tasks," *arXiv preprint arXiv:1710.01813*, 2017.

[14] K. Shiarlis, M. Wulfmeier, S. Salter, S. Whiteson, and I. Posner, "Taco: Learning task decomposition via temporal alignment for control," *arXiv preprint arXiv:1803.01840*, 2018.

[15] D. Wilbers, R. Lioutikov, and J. Peters, "Context-driven movement primitive adaptation," in *International Conference on Robotics and Automation (ICRA)*, 2017.

[16] J. Kober, A. Wilhelm, E. Oztop, and J. Peters, "Reinforcement learning to adjust parametrized motor primitives to new situations," *Autonomous Robots*, 2012.

[17] G. Konidaris, S. Kuindersma, R. Grupen, and A. Barto, "Robot learning from demonstration by constructing skill trees," *The International Journal of Robotics Research*, pp. 360–375, 2012.

[18] A. Nguyen, D. Kanoulas, L. Muratore, D. G. Caldwell, and N. G. Tsagarakis, "Translating videos to commands for robotic manipulation with deep recurrent neural networks," *arXiv:1710.00290*, 2017.

[19] E. E. Aksoy, E. Ovchinnikova, A. Orhan, Y. Yang, and T. Asfour, "Unsupervised linking of visual features to textual descriptions in long manipulation activities," *Robotics and Automation Letters*, 2017.

[20] F. Meier, E. Theodorou, F. Stulp, and S. Schaal, "Movement segmentation using a primitive library," in *Intelligent Robots and Systems (IROS), 2011 IEEE/RSJ International Conference on*. IEEE, 2011, pp. 3407–3412.

[21] P. Pastor, M. Kalakrishnan, L. Righetti, and S. Schaal, "Towards associative skill memories," in *Humanoid Robots (Humanoids), 2012 12th IEEE-RAS International Conference on*. IEEE, 2012, pp. 309–315.

[22] Y. Duan, M. Andrychowicz, B. Stadie, J. Ho, J. Schneider, I. Sutskever, P. Abbeel, and W. Zaremba, "One-shot imitation learning," *Neural Information Processing Systems (NIPS)*, 2017.

[23] D.-A. Huang, S. Nair, D. Xu, Y. Zhu, A. Garg, L. Fei-Fei, S. Savarese, and J. C. Niebles, "Neural task graphs: Generalizing to unseen tasks from a single video demonstration," *arXiv preprint arXiv:1807.03480*, 2018.

[24] J. Butterfield, S. Osentoski, G. Jay, and O. C. Jenkins, "Learning from demonstration using a multi-valued function regressor for time-series data," in *Humanoid Robots (Humanoids), 2010 10th IEEE-RAS International Conference on*. IEEE, 2010, pp. 328–333.

[25] S. Niekum, S. Chitta, A. G. Barto, B. Marthi, and S. Osentoski, "Incremental semantically grounded learning from demonstration." in *Robotics: Science and Systems*, vol. 9. Berlin, Germany, 2013.

[26] K. Hausman, Y. Chebotar, S. Schaal, G. Sukhatme, and J. J. Lim, "Multi-modal imitation learning from unstructured demonstrations using generative adversarial nets," in *Advances in Neural Information Processing Systems*, 2017, pp. 1235–1245.

[27] D. H. Grollman and O. C. Jenkins, "Incremental learning of subtasks from unsegmented demonstration," in *Intelligent Robots and Systems (IROS), 2010 IEEE/RSJ International Conference on*. IEEE, 2010, pp. 261–266.

[28] A. Murali, A. Garg, S. Krishnan, F. T. Pokorny, P. Abbeel, T. Darrell, and K. Goldberg, "Tsc-dl: Unsupervised trajectory segmentation of multi-modal surgical demonstrations with deep learning," in *Robotics and Automation (ICRA), 2016 IEEE International Conference on*. IEEE, 2016.

[29] S. Manschitz, J. Kober, M. Gienger, and J. Peters, "Probabilistic progress prediction and sequencing of concurrent movement primitives," in *Intelligent Robots and Systems (IROS), 2015 IEEE/RSJ International Conference on*. IEEE, 2015, pp. 449–455.

[30] C. Finn, P. Abbeel, and S. Levine, "Model-agnostic meta-learning for fast adaptation of deep networks," *International Conference on Machine Learning (ICML)*, 2017.

[31] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," *International Conference on Learning Representations (ICLR)*, 2015.

[32] P. Ramachandran, B. Zoph, and Q. V. Le, "Swish: a self-gated activation function," *arXiv preprint arXiv:1710.05941*, 2017.

[33] J. L. Ba, J. R. Kiros, and G. E. Hinton, "Layer normalization," *arXiv:1607.06450*, 2016.

[34] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: a simple way to prevent neural networks from overfitting," *The Journal of Machine Learning Research*, vol. 15, no. 1, pp. 1929–1958, 2014.

[35] E. Todorov, T. Erez, and Y. Tassa, "Mujoco: A physics engine for model-based control," in *International Conference on Intelligent Robots and Systems (IROS)*, 2012.

[36] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.

[37] S. Ross, G. J. Gordon, and D. Bagnell, "A reduction of imitation learning and structured prediction to no-regret online learning." in *AISTATS*, 2011.