PLATO: Predicting Latent Affordances Through Object-Centric Play

Suneel Belkhale, Dorsa Sadigh Stanford University

Abstract: Constructing a diverse repertoire of manipulation skills in a scalable fashion remains an unsolved challenge in robotics. One way to address this challenge is with unstructured human play, where humans operate freely in an environment to reach unspecified goals. Play is a simple and cheap method for collecting diverse user demonstrations with broad state and goal coverage over an environment. Due to this diverse coverage, existing approaches for learning from play are more robust to online policy deviations from the offline data distribution. However, these methods often struggle to learn under scene variation and on challenging manipulation primitives, due in part to improperly associating complex behaviors to the scene changes they induce. Our insight is that an objectcentric view of play data can help link human behaviors and the resulting changes in the environment, and thus improve multi-task policy learning. In this work, we construct a latent space to model object affordances – properties of an object that define its uses – in the environment, and then learn a policy to achieve the desired affordances. By modeling and predicting the desired affordance across variable horizon tasks, our method, Predicting Latent Affordances Through Object-Centric Play (PLATO), outperforms existing methods on complex manipulation tasks in both 2D and 3D object manipulation simulation and real world environments for diverse types of interactions. Videos can be found on our website.

Keywords: Human Play Data, Object Affordance Learning, Imitation Learning

1 Introduction

The field of robotics has seen tremendous progress in solving manipulation tasks, but learning a general multi-task policy remains an open challenge. Imitation learning methods are sample-efficient at replicating demonstrated behaviors, but are often presented with structured, predefined tasks and therefore struggle to generalize outside of the data distribution [1, 2]. Rather than using predefined task demonstrations, recent work has shown that learning from *play* data – an unstructured form of demonstration without predefined goals – can lead to policies that are more robust to online deviations [3]. Play data is easy to collect at scale since it requires no task specification or manual resetting, and play can have broad data coverage over the set of object interactions necessary for performing a variety of tasks.

Existing approaches for learning from play sample short horizon-length windows from play data to learn goal-conditioned imitation policies in an offline fashion [3, 4]. However, not all facets of the demonstrator's behavior are captured by the goal alone; thus prior work learns a latent space to model the variation in human behaviors. Such a latent space captures a representation of "plans," for example plans representing reaching or grasping motions, for a given goal state during play. These latent plans can then help inform robot policies at test time [3].

These approaches make several restrictive assumptions. Firstly, they assume that given a fixed short-horizon window of play, the agent's goal is the environment state at the end of the window. This assumption however is not always true: since the desired robot state is not available at test time, the goal is chosen as the environment state a few seconds in the future. Since the goal might be close or the same as the initial state, it can be uninformative for the policy, nor does it necessarily represent the human's true goal. Crediting behavior to an incorrect environment goal can obscure why a human chose the actions during that window, and thus hinder policy learning. For example,

if we want to grasp and move a block on a table but our goal is sampled after the robot has initiated motion but before initiating the grasp, the sampled goal will have no change in the object state and thus gives us no information on the true goal that the user had in mind. Secondly, by randomly sampling windows from play data, the learned latent space will be forced to model all sequences of behaviors *equally*, even though many sequences will be less critical to achieving the desired environment goal. These restrictive assumptions leads to suboptimal behaviors when increasing the complexity and variability of tasks, e.g., tasks with varying horizons, when learning from play.

Instead of defining goals and plans based on arbitrary horizon lengths, we posit that humans often define goals and plan in terms of interactions with objects: rather than planning over the individual joint motions required to grasp and open a door handle, we think about turning the handle and then opening the door. Our key insight is that viewing offline play data as diverse object interactions enables better modeling of both a human's goals and the behaviors that can achieve these goals. Rather than learning to represent fixed short-horizon robot trajectories, or *plans*, we bias the latent space towards learning demonstrated *object affordances* that accomplish tasks in object-space. Imitation policies can then condition on these affordances directly to gain insight into human's behaviors.

We propose an algorithm, PLATO – Predicting Latent Affordances Through Object-Centric Play – that automatically segments play into a series of object interactions using proprioceptive information, and then learns a latent affordance space over object interactions. Simultaneously, it learns to imitate human actions over variable horizon interactions, conditioned on the latent affordance and goal. By considering object interactions in play and correctly attributing goals to robot behaviors, PLATO builds a robust mapping between the true goals, desired object affordances, and actions on the robot over varying horizons. This leads to PLATO significantly outperforming prior methods especially when increasing the complexity and variability of play data. Our contributions are:

- 1. We have developed an object-centric paradigm for learning from unstructured human play data, which views play as sequential, unlabeled interactions with objects.
- 2. We have developed a new algorithm, PLATO, that extracts and leverages these interactions to model diverse object affordances from play data and learn a robust policy.
- 3. We have tested PLATO on a number of 2D and 3D manipulation environments in simulation and the real world, including diverse objects such as blocks, mugs, cabinets, and drawers, with broad coverage over possible tasks and object affordances. We demonstrate that PLATO substantially outperforms state-of-the-art learning from play baselines in these complex manipulation tasks.

2 Related Work

In this section, we will discuss prior work in goal-conditioned imitation learning, learning from play, and object-centric policy learning. Our work brings an object-centric perspective on learning from play to learn effective imitation policies that generalize across different manipulation tasks.

Imitation Learning. Imitation Learning is a common method for learning robot policies from human demonstrations, where a policy learns to mimic human actions [1]. These methods often struggle to generalize to new environments and to learn from multi-task data [5, 6, 2]. To improve *generalization*, Ho, et al. introduced generative adversarial IL to learn from imitation data while matching the expert policy distribution [7]. Recently, implicit imitation learning policies using energy-based models have also been shown to improve generalization as compared to explicit policy models [8]. To enable *multi-task* learning, one can condition the policy on goal states, either explicitly labelled or inferred via hindsight experience replay [9, 10, 11]. Other methods have leveraged meta-learning for multi-task imitation learning to enable better multi-task performance and one-shot generalization [12]. A key limitation of these works is the dependency on demonstration quality and quantity, for example the state-action coverage within the dataset for each task [13]. Our work utilizes an approach built on top of goal-conditioned imitation learning to learn *multi-task* policies. However, we use play data as opposed to expert demonstrations, which enables a much broader state-action coverage to learn *generalizable* and robust policies.

Learning from Play. Human play data is defined as unstructured and unsupervised human interaction with an environment, as a means to let the user guide the data collection process, and provide broader coverage of the task-relevant state and goal space. Because humans can freely choose how to interact with the environment, play data is easy and cheap to collect without the need for prior task specification or environment resets. These factors make play data a rich source for skill learning; as

a result, imitation policies trained on play data are more robust to deviations from the expert trajectories than those trained on single-task demonstrations [3]. The state of the art method, Play-LMP, learns latent "plans" over short, fixed horizon trajectories from play data [3]. This approach can also be used to ground language with play data using pretrained language models [14]. These learned short-horizon skills have also been chained together to accomplish long horizon tasks using motion planning inspired techniques [15]. Additionally, using small variations in the horizon length of play windows has been shown to improve test time performance [4]. These methods benefit from the broad coverage of play data, but suffer from distribution shift at test time due to incorrectly inferred goals during training. Therefore, our approach also imitates play, but we take an object-centric view: by learning from diverse interactions with the environment, we can more accurately infer the human's goals during play and thus obtain a better state-action-goal distribution.

Object-Centric Modeling and Policy Learning. There is strong evidence that humans have neural pathways specific to object recognition, dynamics understanding, and scene segmentation [16]. Inspired by humans, this notion of object centrism has been applied in many recent works in object manipulation for better policy generalization and sample efficiency. Formalisms like object-oriented MDPs have been introduced to better model *effects* of agent behaviors as functions of individual objects and their affordances [17]. Object affordances have also been learned from labeled interaction examples in multi-task environments and were shown to benefit policy learning downstream [18, 19]. Planning relative to object reference frames and primitives can enable plans to generalize to changing environments [20, 21]. Similarly, focusing on objects during policy learning can allow behaviors to generalize to novel scenarios [22, 23]. A recent work leveraged labelled object motions in play data to predict grasp points, and showed RL becomes more sample efficient post-grasping [24]. We employ object centrism to guide learning from play data, specifically by extracting affordances through unsupervised temporal segmentation over object interactions.

3 Predicting Latent Affordances Through Object-Centric Play

A play dataset D_{play} consists of N varying length episodes of undirected, human generated state-action trajectories $T_j, \ \forall j \in \{1,\dots,N\}$, where $T_j = \{s_1,a_1,\dots,s_{L_j}\}$ for length L_j . The state feature space S may be learned or predefined, and consists of a proprioceptive state space S^r (robot state space) and an environmental state space S^o (object state space) such that $S^r \oplus S^o = S$. The goals that generated these trajectories are not included in the dataset, and are defined as $o_g \in G$, for goal space $G \subseteq S^o$. We use o and s^o interchangeably to refer to environmental state, and likewise for r and s^r to refer to the robot state. We assume that goals only consist of the environmental state, since access to proprioceptive state goals at test time is unrealistic as it requires the robot having access to a policy for achieving the given environmental goal.

Given access to this play dataset D_{play} , a new initial state $s \in S$, and an object goal state o_g , our problem is to learn a robot policy π to achieve the desired goal state. In prior work, sampled trajectories of play τ consist of a contiguous fixed-horizon segment from an episode T_j . For simplicity we denote sampling these segments from play as $\tau \sim D_{\text{play}}$, where the length of the segment is the fixed horizon H. We propose viewing play sequences from a bird's eye view, namely object interactions, instead of from myopic and fixed 1-2 second windows used in prior work. From this perspective, we hypothesize human play is just a series of *environment interactions* induced by a robot's actions. If we can somehow detect where interactions begin and end, we can learn to relate interaction and pre-interaction robot behaviors to the state of the world post-interaction. First, we discuss how we detect interactions in the environment in offline play by leveraging proprioceptive cues. Next, we formalize how we properly choose and associate goals with robot behaviors.

Segmenting Play into Interaction Phases: When we interact with an environment, we change its state through our own actions, whether it be through direct or indirect contact. In this work, we focus on single-object interactions, but we emphasize that interactions can be defined even over multiple objects that are being influenced by the robot behavior (e.g., in tool-mediated manipulation, interaction is defined between the tool and the environment). We break down an interaction into the following phases:

- 1. **Pre-interaction:** This phase usually involves orienting the robot to interact with an object, e.g., reaching the purple block in Fig. 1.
- 2. **Interaction:** This phase involves joint and interdependent motions between the robot and the object(s), e.g., pulling the purple block in Fig. 1.

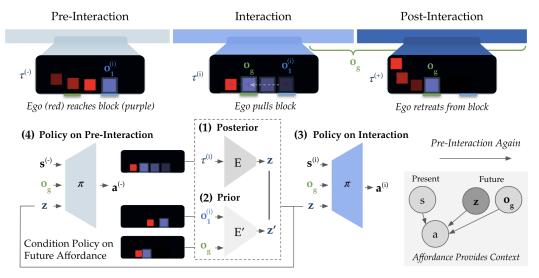


Figure 1: Our method shown on pre-interaction (purple), interaction (blue), and post-interaction (green) periods. (1) Posterior E encodes the interaction sequence $\tau^{(i)}$ into affordance z. (2) Prior E' encodes object start and goal states $o_1^{(i)}$ and o_g to predict z. o_g is sampled from post-interaction. (3) Policy trained to output actions on interaction period conditioned on affordance, and simultaneously (4) to output actions from pre-interaction period conditioned on the "future" affordance. The assumed causal structure is shown in the lower right: PLATO claims the robot behavior can be explained with knowledge of the long term goal and future affordance z: the policy reasons about its desired affordance and goal to determine its actions.

3. **Post-interaction:** This phase involves separation from the object and any downstream effects on the object, e.g., the purple block coming to rest after releasing it in Fig. 1. For repetitive interactions, this is the same as the next pre-interaction period.

To account for these periods, we detect the robot's *influence* over the environment, e.g., grasping a door handle or lifting a mug, to automatically segment play into the periods above. For detecting single-object interactions, contact signals are readily available with modern robots. The contact signal is smoothed and then chunks of play with contiguous contact are labelled as interaction periods. Several ideas for extending PLATO to handle multi-object interaction detection are discussed in detail in Appendix E. Since play consists of many back-to-back interactions, the post interaction phase of one segment, e.g., retreating from a block and moving towards a cup, is simply the preinteraction for the next segment, e.g., lifting the cup. We denote pre, interaction, and post windows with superscripts ⁽⁻⁾, ⁽ⁱ⁾, and ⁽⁺⁾, respectively.

Sampling Task Relevant Goals Post-Interaction: By segmenting play into interaction phases, we can more accurately sample an accurate goal environment state. Critically, goals *result* from an interaction by definition. At the end of an interaction, the environment state will very likely have changed. Therefore, rather than arbitrarily labeling the goal as the last state of a 1-2 second window, as in prior work, we can label the goal as any state from the end of the interaction through post-interaction (see o_g in Figure 1). This includes any downstream effects on the object from the interaction (e.g., gravity or inertia). For example, if we slide a block along a table, a valid o_g is any block state before the block stops sliding. With an informative goal sampled, we can proceed to learn a goal-conditioned policy.

Extracting Affordances from Interaction for Robust Policy Learning: We define every interaction between the robot and the environment as exploiting some *affordance* on an object in the scene. Affordances are properties of objects that define how they can be used (e.g., a block being grasped, a door knob being turned, or a drawer being opened). Our insight is that learning these affordances (what happens to the object) instead of plans (what happens to the robot) from play will lead to a much simpler and more robust task representation that can operate over varying horizons, and thus will yield much better policies at test time. This paradigm empowers the policy to *reason* about the environment: given access to an affordance (e.g., the door knob being turned) and the goal (e.g., opened door), the policy should be able to work backwards to infer the behavior to exploit that affordance (e.g., reach the knob and rotate the gripper to turn it). This is in contrast to prior work that

relies on randomly selected, short, fixed horizon windows to learn latent representations of plans such plans fail to capture varying horizon tasks and overly depend on the robot state, leading to generalization issues at test time [3].

PLATO Design: The PLATO architecture is shown in Fig. 1. At a high level, PLATO learns to model each interaction trajectory $\tau^{(i)}$ as a latent affordance z ((1) **Posterior** in Fig. 1). This latent affordance z explains the actions both leading up to (pre-interaction) and during the interaction (see causal structure in Fig. 1). For example, knowing that we want to push a block (goal) and how we want to push it (affordance) allows the policy to infer that it should servo to the correct side of the block first. Therefore, the latent space is learned end-to-end with the policy π , which decodes latent affordances conditioned on the current state and object goal states to reproduce both pre-interaction and interaction robot actions ((3) Policy on Interaction and (4) Policy on Pre-Interaction in Fig. 1). By not learning to encode pre-interaction sequences, the critical assumption we make here is that variations in the affordances (interaction phase) are more relevant to the task and thus more important to capture in our latent space than variations across any random behavior sequence (pre-interaction).

The policy action-reconstruction objective alone would encourage z to model robot behaviors (plans), leading to a myopic view of the task. In order to force the latent space to focus on object affordances and thus be more robust at test time, we simultaneously learn a prior on the affordance distribution conditioned on just the start and goal object states ((2) Prior in Fig. 1), trained to regularize the posterior latent affordance z. The posterior on z considers the full window $\tau^{(i)}$, while the prior sees just the start and goal object states, and so the prior helps shape the latent space to encode the affordance in $\tau^{(i)}$ rather than just the action information.

Algorithm 1 PLATO Training

```
1: Given: H^{(i)}, H^{(-)}, play data D_{\text{play}}, interaction criteria f^{(i)}, 2: D_{\text{play}}^{(-)}, D_{\text{play}}^{(i)}, D_{\text{play}}^{(+)} = f^{(i)}(D_{\text{play}})
3: Initialize E \cap E'
                                                                                                                                                                                            ▷ Split into interactions
  3: Initialize E, E', \pi
  4: while not converged do
 5: \tau^{(-)}, \tau^{(i)}, \tau^{(+)} \sim D_{\text{play}}^{(-)}, D_{\text{play}}^{(i)}, D_{\text{play}}^{(+)}
6: Sample o_g \sim \{o_t^{(+)}\}
7: p(z) \leftarrow E(\tau^{(i)})
8: p(z') \leftarrow E'(o_1^{(i)}, o_g)
                                                                                                                                                                  ▶ Posterior Affordance Distribution
                                                                                                                                                                           ▶ Prior Affordance Distribution
9: z \sim p(z)

10: \tilde{a}_{1:H^{(i)}}^{(i)} \leftarrow \pi(s_{1:H^{(i)}}^{(i)}, o_g, z)

11: \tilde{a}_{1:H^{(-)}}^{(-)} \leftarrow \pi(s_{1:H^{(-)}}^{(-)}, o_g, z)
                                                                                                                                                                                              ▶ Policy in Interaction
                                                                                                                                                                                      ⊳ Policy in Pre-Interaction
               Compute \mathcal{L}_{PLATO} with Eq. (1) and update \pi, E, E'.
```

The training procedure is outlined in Alg. 1. After sampling windows from each interaction phase, $\tau^{(-)}$, $\tau^{(i)}$, and $\tau^{(+)}$ (Line 5), as well as a long term goal o_a (Line 6), PLATO encodes the interaction into a posterior and prior affordance distribution (Lines 7-8). Next an affordance z is sampled from the posterior using the reparameterization trick, and z is then used to decode actions during interaction (Line 10) and pre-interaction (Line 11). See Appendix B.1 for more discussion of this training procedure, including its computational efficiency. During training, we utilize action reconstruction losses over both the pre-interaction region and the interaction region windows to train all networks end-to-end. We utilize a KL divergence term to both train the affordance prior network and regularize the affordance posterior network.

$$\mathcal{L}_{\text{PLATO}} = -\log(\pi(a_{1:H}^{(i)}|s_{1:H}^{(i)}, o_g, z)) - \alpha\log(\pi(a_{1:H}^{(-)}|s_{1:H}^{(-)}, o_g, z)) + \beta \text{ KL}(p(z) \parallel p(z'))$$
 (1)

Here, α controls the policy focus on reconstructing pre-interaction behaviors, which usually is set to 1. This procedure of learning a posterior and prior network and conditioning the policy on the latent context resembles existing work [3]. The key differences are i) our goals are sampled based on interactions with objects, ii) our latent representations capture object-centric affordances z by intelligently shaping the learning objective, which can lead to more effective and generalizable policies, and iii) the policy reasons over varying and longer horizon sequences by conditioning on the same future desired affordance z and long term goal o_g throughout both pre-interaction and interaction.

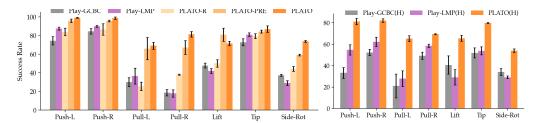


Figure 2: Block2D Success Rates, trained over 3 random seeds and evaluated on various primitives for PLATO and baselines Play-LMP and Play-GCBC. **Left:** Scripted play data, with ablations PLATO-PRE and PLATO-R. **Right:** Human play data. PLATO substantially outperforms baselines on both scripted and human play data.

Test Time: At test time, the affordance posterior network E cannot be used to output plans z, since it assumes access to a trajectory. Similar to prior work [3], the prior network E' is used to propose z' at test time, using only the current object state o_t and the desired goal object state (o_g) . Unlike prior work, the prior network only depends on object states, and will be robust to different robot starting states. Therefore instead of having to predict the exact robot behavior (i.e., plan) to achieve a goal, PLATO predicts an affordance at test time and empowers the policy to exploit this affordance. The policy conditions on z' and goal o_g to produce actions at the current state.

Addressing Challenges of Prior Work: By choosing goals o_g resulting from interactions with objects, PLATO better reflects the true demonstrator goal, and thus can reduce the credit assignment problem found with fixed-horizon methods. PLATO handles variable horizon sequences by explicitly training the policy on pre-interaction and interaction sequences together, conditioned on the affordance and long term goal. By biasing the latent space to model *affordances* during interaction, rather than just any random sequence of play, PLATO learns the task relevant behaviors and variations therein to aid in generalization. Our results in Sec. 4 demonstrate the effectiveness of PLATO compared to prior work on a wide range of complex tasks. See Appendix A for further discussion.

4 Experiments

In this section, we evaluate our approach extensively across three single-object manipulation environments (Block2D, Block3D-Platforms, Mug3D-Platforms), one multi-object scene (Playroom3D), and one real scene (Block-Real), across diverse tasks like pushing, lifting, and rotating. These environments, shown in Fig. 3, enable a wide variety of objects and possible affordances. We collect scripted play data in all environments as well as human play data for Block2D, and train each method until convergence. See Appendix C for environment, task, data collection, and training details, and Appendix B.2 for method implementation details.

Baselines: We compare PLATO against the two state-of-the-art methods for Learning from Play, Play-GCBC (Goal-Conditioned BC) and Play-LMP (Latent Motor Plans) [3]. We also implement two variants of our method, PLATO-PRE, which encodes both the interaction *and* the pre-interaction periods into the latent space (adding $\tau^{(-)}$ as an input to E in step (1) in Fig. 1), and PLATO-R, which replaces the object-centric prior with the prior from Play-LMP (adding the initial robot state as input to E' in Fig. 1). PLATO-PRE results show the sufficiency of affordances as a latent representation, while PLATO-R results show how object-centrism benefits the learned latent space and the policy.

Block2D Results: Evaluation results for PLATO for both scripted play data and human play data are shown in Fig. 2, along with results on PLATO-PRE and PLATO-R. On scripted play data, PLATO is able to substantially outperform the baselines on every task. Learning affordances from interactions enables PLATO to model much more complex action sequences, such as those involving the tether action, with high accuracy, even across variations in object dimensions, masses, and initial conditions. On human play data, we again find that PLATO outperforms Play-GCBC and Play-LMP on all of the tasks, emphasizing the ability of our method to scale to human generated data. Interestingly, performance for all methods is worse on the human generated data than scripted data. We attribute this to human data containing many sub-optimal trajectories due to the challenges of teleoperation.

PLATO-PRE, which encodes the pre-interaction period, performs slightly worse but similar to PLATO, validating our hypothesis that interaction trajectories (PLATO) contain sufficient information about the task when compared to also including pre-interaction trajectories in the latent space

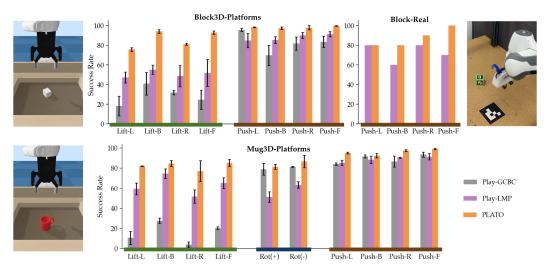


Figure 3: 3D Environment Success Rates, trained over 3 random seeds and evaluated on various primitives for PLATO and baselines Play-LMP and Play-GCBC. **Top Left:** Block3D-Platforms. **Bottom:** Mug3D-Platforms. PLATO substantially outperforms baselines in 3D manipulation environments for pushing, rotating, and lifting tasks. **Top Right:** Block-Real. PLATO trained only in simulation generalizes to real world pushing tasks.

(PLATO-PRE). We find that by adding robot state information to the prior (PLATO-R), performance suffers on most tasks. We hypothesize that this phenomenon is caused by the prior relying too much on the initial state of the robot, and not enough on that of the object. As a result, the latent space will not generalize at test time when the agent inevitably diverges from the offline state distribution. Overall, we see that framing play data through the lens of object interactions (both PLATO and PLATO-PRE) results in much better policy learning than prior state-of-the-art methods.

Block3D-Platform & Mug3DPlatform Results: Fig. 3 shows the results of evaluating PLATO on the harder Platform tasks. Interestingly, Both Play-GCBC and Play-LMP do well on the pushing tasks in this setting, but do very poorly on the more complex lifting tasks (and rotate tasks for Mug3D). By modeling an affordance space and properly relating these affordances to goal environment states, our method is capable of recreating diverse types of varying horizon behaviors from play, and unlike Play-LMP, our method scales smoothly as the number of tasks increases.

Playroom3D Results: Fig. 4 shows the results on the challenging Playroom3D environment with the drawer, cabinet, and the block, with especially diverse affordances and varying horizon tasks. Once again, PLATO outperforms the baselines on every task. While neither baseline can perform the complex Cabinet-Close primitive, PLATO achieves 100% success. Notably, Play-GCBC performs better than Play-LMP on several tasks, suggesting that the prior network and policy might be out of distribution for the test time states and goals on these complex primitives.

Block-Real Results: Fig. 3 shows the results of evaluating PLATO on pushing tasks for a real robot setup with *no additional real world data* (see Appendix C for details). Play-LMP and PLATO both get near perfect success in simulation since these pushing tasks are less complex, but PLATO generalizes better across the gap between real and simulated object dynamics.

Overall, PLATO achieves substantially higher success rates than the baselines on a wide variety of tasks and object properties in simulation and real environments, along with lower variance across random seeds. This trend is even more apparent when we go beyond simple pushing tasks and consider more complex object interactions such as opening and closing doors, drawers, cabinets, lifting, and others. Due to its object-centric view of play, PLATO is able to extract and exploit diverse types of affordances in the environment and can handle varying horizon tasks. See Appendix D for additional experiments including more tasks per scene and a longer discussion of results.

5 Conclusion

Summary: In this work, we introduced an object-centric paradigm for learning from play data involving segmenting play into a series of object interactions. Prior state-of-the-art methods for

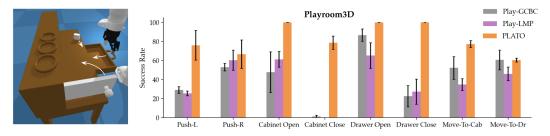


Figure 4: Playroom3D 3D Environment Success Rates. The difference between PLATO and baselines is even greater for more complex and varying horizon tasks, like opening and closing drawers and cabinets. See Appendix D.1 for additional tasks in this environment involving object retrieval and button pressing.

learning from play suffer from credit assignment issues that stem from the short and fixed horizon of sampled trajectories. Our method, PLATO, addresses these credit assignment issues by choosing goals that involve meaningful changes in object state. PLATO learns a latent affordance space to model these interactions and their variations, and simultaneously learns to predict these latent affordances from the goal object state. These latent affordances help to inform the robot behavior across varying horizon tasks. Through our extensive experiments in both 2D, 3D, and real world environments spanning a wide variety of object manipulation tasks, we show that PLATO substantially outperforms prior methods on both scripted and human play data.

Limitations and Future Work: Our work introduces a paradigm of learning from variable horizon object interactions, and our method achieves substantially better performance across a variety of *single-object* manipulation tasks. In future work, we intend to expand our notion of interaction to encompass even *multi-object* interactions in play (e.g., tool use). When using tools, the robot might have second and third order effects on the environment that we could model. We believe the interaction paradigm introduced in this work is an important first step to reasoning about these higher order effects. For both the single and multi-object settings, future work might leverage notions of action information density or find bottleneck states to automatically segment interactions.

As shown in Section 4, data collection methods for play greatly affect final policy performance. We primarily evaluate with scripted play, but we hope to collect large human play datasets in future work. Compared to scripted play, we hypothesize that human play consists of much more behavior variability, yielding significant plan and affordance variability for a given start and goal state. Isolating affordances and interactions helps manage this variability, and PLATO is still able to perform well on all the tasks as a result. However, a future direction would be to study how exactly human play data differs from machine generated play data in order to develop more robust methods.

To add, our real robot experiments demonstrate PLATO trained in simulation can generalize to real world dynamics for pushing tasks without any data from the real robot. We hope to collect play data directly on the robot in future work to learn even more complex tasks in the real world. Furthermore, our method makes use of object state information, which may not always be easily available in practice. However, we claim this method can readily handle images with several simple changes to the learned prior, described in Appendix E.

We present thorough discussions of these limitations as well as potential solutions in Appendix E.

Acknowledgements

This work was funded by JP Morgan, the Office of Naval Research, DARPA YFA, and NSF Award Numbers 1849952, 1941722, and 2006388.

References

- [1] D. A. Pomerleau. Alvinn: An autonomous land vehicle in a neural network. In D. Touretzky, editor, *Advances in Neural Information Processing Systems*, volume 1. Morgan-Kaufmann, 1989. URL https://proceedings.neurips.cc/paper/1988/file/812b4ba 287f5ee0bc9d43bbf5bbe87fb-Paper.pdf.
- [2] S. Ross and D. Bagnell. Efficient reductions for imitation learning. In Y. W. Teh and M. Titterington, editors, *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, volume 9 of *Proceedings of Machine Learning Research*, pages 661–668, Chia Laguna Resort, Sardinia, Italy, 13–15 May 2010. PMLR. URL https://proceedings.mlr.press/v9/ross10a.html.
- [3] C. Lynch, M. Khansari, T. Xiao, V. Kumar, J. Tompson, S. Levine, and P. Sermanet. Learning latent plans from play. In *Conference on Robot Learning*, pages 1113–1132. PMLR, 2020.
- [4] A. Gupta, V. Kumar, C. Lynch, S. Levine, and K. Hausman. Relay policy learning: Solving long-horizon tasks via imitation and reinforcement learning. In L. P. Kaelbling, D. Kragic, and K. Sugiura, editors, 3rd Annual Conference on Robot Learning, CoRL 2019, Osaka, Japan, October 30 November 1, 2019, Proceedings, volume 100 of Proceedings of Machine Learning Research, pages 1025–1037. PMLR, 2019. URL http://proceedings.mlr.press/v100/gupta20a.html.
- [5] S. Schaal. Is imitation learning the route to humanoid robots? Trends in Cognitive Sciences, 3(6):233-242, 1999. ISSN 1364-6613. doi:https://doi.org/10.1016/S1364-6613(99)01327-3. URL https://www.sciencedirect.com/science/article/pii/S1364661 399013273.
- [6] B. D. Argall, S. Chernova, M. Veloso, and B. Browning. A survey of robot learning from demonstration. *Robotics and Autonomous Systems*, 57(5):469–483, 2009. ISSN 0921-8890. doi:https://doi.org/10.1016/j.robot.2008.10.024. URL https://www.sciencedirect. com/science/article/pii/S0921889008001772.
- [7] J. Ho and S. Ermon. Generative adversarial imitation learning. In D. D. Lee, M. Sugiyama, U. von Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pages 4565–4573, 2016. URL https://proceedings.neurips.cc/paper/2016/hash/cc7e2b878868cbae992d1fb743995d8f-Abstract.html.
- [8] P. Florence, C. Lynch, A. Zeng, O. A. Ramirez, A. Wahid, L. Downs, A. Wong, J. Lee, I. Mordatch, and J. Tompson. Implicit behavioral cloning. In A. Faust, D. Hsu, and G. Neumann, editors, *Conference on Robot Learning*, 8-11 November 2021, London, UK, volume 164 of Proceedings of Machine Learning Research, pages 158–168. PMLR, 2021. URL https://proceedings.mlr.press/v164/florence22a.html.
- [9] F. Codevilla, M. Müller, A. M. López, V. Koltun, and A. Dosovitskiy. End-to-end driving via conditional imitation learning. In 2018 IEEE International Conference on Robotics and Automation, ICRA 2018, Brisbane, Australia, May 21-25, 2018, pages 1–9. IEEE, 2018. doi: 10.1109/ICRA.2018.8460487. URL https://doi.org/10.1109/ICRA.2018.8460487.
- [10] Y. Ding, C. Florensa, P. Abbeel, and M. Phielipp. Goal-conditioned imitation learning. In H. M. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. B. Fox, and R. Garnett, editors, Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada, pages 15298–15309, 2019. URL https://proceedings.neurips.cc/paper/2019/hash/c8d3a760ebab631565f8509d84b3b3f1-Abstract.html.
- [11] M. Andrychowicz, D. Crow, A. Ray, J. Schneider, R. Fong, P. Welinder, B. McGrew, J. Tobin, P. Abbeel, and W. Zaremba. Hindsight experience replay. In I. Guyon, U. von Luxburg, S. Bengio, H. M. Wallach, R. Fergus, S. V. N. Vishwanathan, and R. Garnett, editors, *Advances*

- in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA, pages 5048–5058, 2017. URL https://proceedings.neurips.cc/paper/2017/hash/453fadbd8a1a3af50a9df4df899537b5-Abstract.html.
- [12] Y. Duan, M. Andrychowicz, B. C. Stadie, J. Ho, J. Schneider, I. Sutskever, P. Abbeel, and W. Zaremba. One-shot imitation learning. In I. Guyon, U. von Luxburg, S. Bengio, H. M. Wallach, R. Fergus, S. V. N. Vishwanathan, and R. Garnett, editors, Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA, pages 1087–1098, 2017. URL https://proceedings.neurips.cc/paper/2017/hash/ba3866600c3540f 67cle9575e213be0a-Abstract.html.
- [13] A. Mandlekar, D. Xu, J. Wong, S. Nasiriany, C. Wang, R. Kulkarni, L. Fei-Fei, S. Savarese, Y. Zhu, and R. Martín-Martín. What matters in learning from offline human demonstrations for robot manipulation. In A. Faust, D. Hsu, and G. Neumann, editors, *Conference on Robot Learning*, 8-11 November 2021, London, UK, volume 164 of Proceedings of Machine Learning Research, pages 1678–1690. PMLR, 2021. URL https://proceedings.mlr.press/v164/mandlekar22a.html.
- [14] C. Lynch and P. Sermanet. Language conditioned imitation learning over unstructured data. In D. A. Shell, M. Toussaint, and M. A. Hsieh, editors, *Robotics: Science and Systems XVII*, *Virtual Event, July 12-16, 2021*, 2021. doi:10.15607/RSS.2021.XVII.047. URL https://doi.org/10.15607/RSS.2021.XVII.047.
- [15] B. Ichter, P. Sermanet, and C. Lynch. Broadly-exploring, local-policy trees for long-horizon task planning. In A. Faust, D. Hsu, and G. Neumann, editors, *Conference on Robot Learning, 8-11 November 2021, London, UK*, volume 164 of *Proceedings of Machine Learning Research*, pages 59–69. PMLR, 2021. URL https://proceedings.mlr.press/v164/ichter22a.html.
- [16] K. Grill-Spector. The neural basis of object perception. Current Opinion in Neurobiology, 13 (2):159-166, 2003. ISSN 0959-4388. doi:https://doi.org/10.1016/S0959-4388(03)00040-0. URL https://www.sciencedirect.com/science/article/pii/S0959438 803000400.
- [17] C. Diuk, A. Cohen, and M. L. Littman. An object-oriented representation for efficient reinforcement learning. In *Proceedings of the 25th international conference on Machine learning*, pages 240–247, 2008.
- [18] T. Nagarajan and K. Grauman. Learning affordance landscapes for interaction exploration in 3d environments. In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, editors, Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020. URL https://proceedings.neurips.cc/paper/2020/hash/15825aee15eb335cc13f9b559f166ee8-Abstract.html.
- [19] A. Khazatsky, A. Nair, D. Jing, and S. Levine. What can I do here? learning new skills by imagining visual affordances. In *IEEE International Conference on Robotics and Automation, ICRA 2021, Xi'an, China, May 30 June 5, 2021*, pages 14291–14297. IEEE, 2021. doi: 10.1109/ICRA48506.2021.9561692. URL https://doi.org/10.1109/ICRA48506.2021.9561692.
- [20] T. Migimatsu and J. Bohg. Object-centric task and motion planning in dynamic environments. *IEEE Robotics and Automation Letters*, 5(2):844–851, Apr 2020. ISSN 2377-3774. doi:10.1 109/lra.2020.2965875. URL http://dx.doi.org/10.1109/lra.2020.2965875.
- [21] M. Dalal, D. Pathak, and R. R. Salakhutdinov. Accelerating robotic reinforcement learning via parameterized action primitives. Advances in Neural Information Processing Systems, 34, 2021.

- [22] C. Devin, P. Abbeel, T. Darrell, and S. Levine. Deep object-centric representations for generalizable robot learning. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 7111–7118, 2018. doi:10.1109/ICRA.2018.8461196.
- [23] M. Zhao, Z. Liu, S. Luan, S. Zhang, D. Precup, and Y. Bengio. A consciousness-inspired planning agent for model-based reinforcement learning. *Neural Information Processing Systems*, abs/2106.02097, 2021. URL https://papers.nips.cc/paper/2021/hash/0c2 15f194276000be6a6df6528067151-Abstract.html.
- [24] J. Borja-Diaz, O. Mees, G. Kalweit, L. Hermann, J. Boedecker, and W. Burgard. Affordance learning from play for sample-efficient policy learning. In 2022 IEEE International Conference on Robotics and Automation, ICRA. IEEE, 2022. URL https://arxiv.org/abs/2203.00352.
- [25] Victor Blomqvist. Pymunk 2d physics engine. URL http://www.pymunk.org/en/la test/index.html.
- [26] T. Kipf, Y. Li, H. Dai, V. Zambaldi, A. Sanchez-Gonzalez, E. Grefenstette, P. Kohli, and P. Battaglia. Compile: Compositional imitation learning and execution. In *International Conference on Machine Learning*, pages 3418–3428. PMLR, 2019.
- [27] R. Ranganath, S. Gerrish, and D. Blei. Black Box Variational Inference. In S. Kaski and J. Corander, editors, *Proceedings of the Seventeenth International Conference on Artificial Intelligence and Statistics*, volume 33 of *Proceedings of Machine Learning Research*, pages 814–822, Reykjavik, Iceland, 22–25 Apr 2014. PMLR. URL http://proceedings.mlr.press/v33/ranganath14.html.
- [28] S. Douglas. Transfer learning from play and language nailing the baseline, Mar 2021. URL https://sholtodouglas.github.io/Learning-from-Play/#what-took-us-so-long.

A Prior Work: Summary and Challenges

Goal-conditioned behavior cloning on play trajectories $\tau \sim D_{\text{play}}$ (Play-GCBC) is quite effective at learning a robust multi-task policy; however, since the policy is only conditioned on the goal, it fails to capture all the degrees of behavior variation that would achieve that goal [3]. Play-LMP addresses this by introducing a latent plan to capture not just *what* goal to reach, but also *how* to reach it. At a high level, Play-LMP learns to represent sampled play trajectories $\tau \sim D_{\text{play}}$ of fixed horizon H as plans in a latent space, denoted by vector z, using variational inference techniques. The method trains a behavior cloning agent to reproduce the humans actions $a_{1:H} \in \tau$ conditioned on the current plan z and the hindsight-labelled goal state o_H .

At a high level, Play-LMP encodes random sequences of environment states and robot actions as latent "plans," which then get passed to a policy that learns to decode these plans at each state into the corresponding human demonstrated action for that time step. They simultaneously learn to predict the latent plan from just the initial and final state in the environment, for use at test time. In practice, the sequences are uniformly sampled 1-2 second chunks from play. Our method, PLATO, samples *interactions* and intelligently learns a latent space from them, enable longer and variable horizon views of sequences of play.

Architecture: To represent play trajectories $\tau \sim D_{\text{play}}$ as plans, Play-LMP uses a posterior encoder E that maps the states $s_{1:H} \in \tau$ to a single latent plan distribution $z \sim \mathcal{N}(\mu_z, \sigma_z^2)$. To regularize the posterior, Play-LMP simultanously learns prior network E' that takes in just the first state s_1 and the goal environment state $s_H^o \in S^o$ to produce the latent plan distribution $z' \sim \mathcal{N}(\mu_{z'}, \sigma_{z'}^2)$. To decode plans into actions at a given state, Play-LMP utilizes a policy π that takes in a plan z, sampled from either the prior distribution (test time) or posterior distribution (train time), along with the current state s_t to predict the action a_t at that step.

E, E', and π are recurrent networks that operate over the fixed time horizon H and are learned end-to-end at training time. This architecture can be viewed as "encoding" trajectories in the environment into a lower dimensional plan z, and then reconstructing these plans into actions with the policy. PLATO similarly learns posterior, prior, and policy networks, but uses different inputs to each network based on sampling interactions and more meaningful goal environment states, as seen in Figure 1.

Training: E, E', and π are trained end-to-end by minimizing the following loss, equivalent to maximizing the evidence-based lower bound (ELBO) of the data likelihood under the posterior network E, given the learned prior network E':

$$\mathcal{L}_{LMP} = -\frac{1}{H} \sum_{t=0}^{H-1} \log(\pi(a_t | s_t, s_H^o, z)) + \beta \text{ KL}(\mathcal{N}(\mu_z, \sigma_z^2) \parallel \mathcal{N}(\mu_{z'}, \sigma_{z'}^2))$$
 (2)

There may be a variety of ways to go from the initial state (s_0) to the final state (s_H^o) ; for example, in order to move a block to the right, we might grab-move the block or push block without grabbing. Play-LMP benefits from encoding these variations in the latent space. PLATO uses a similar variational loss as Play-LMP, but instead of reconstructing actions from the short, fixed horizon τ , reconstructs the interaction and pre-interaction trajectories $\tau^{(-)}$ and $\tau^{(i)}$ conditioned on the affordance extracted from $\tau^{(i)}$ (see Eq. 1).

Test Time: At test time, the plan posterior network E cannot be used to output plans z, since it assumes access to the full trajectory. Therefore, the plan prior network E' is used to propose z' at test time, using only the initial (current) state s_t and the desired goal environment state (s_H^o) . The policy then uses the resulting distribution over z' to predict actions online.

A.1 Play-LMP Strengths

Play-LMP is able to interpret the effects of short interactions with the environment. In doing so, it learns to propose plan distributions p(z') for achieving a wide variety of goals using just the learned prior E'. By capturing all the possible variations of going between each pair of start states s_0 and goal environment states s_H^o during training, the Play-LMP policy is capable of executing a wide variety of behaviors at test time. As compared to Play-GCBC on sampled play windows, Play-LMP performs better on a wide range of manipulation tasks. Play-LMP and Play-GCBC are also

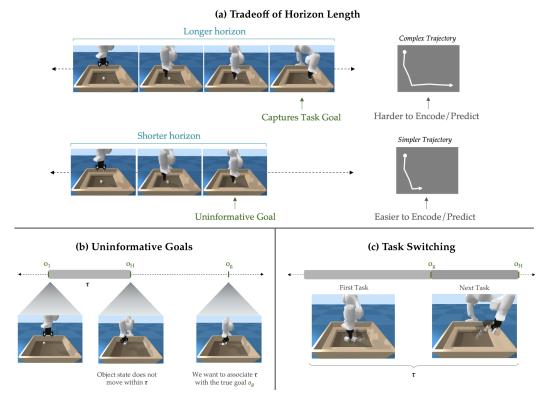


Figure 5: Challenges of short, fixed horizons, visualized for an example pushing task. (a) Choosing the horizon act is a balancing act between capturing task and goal information in the plan, and being able to predict plans (longer horizon means more complexity to capture) from just the initial and goal states. (b) If the horizon length is shorter than the current task, the goal environment state might be the same as the current state (if object state doesn't change within tau), and thus we will not properly link behaviors with the goals that induced these behaviors. (c) If the horizon is too long, We might capture multiple tasks within tau, and thus improperly assign the goal for the first task as being that of the next task. Our method, PLATO, addresses these challenges by biasing the latent space to learn from *object interactions* within play.

shown be robust to minor perturbations in the environment, which can be attributed to the broad state-action-goal distributions found in human play data.

A.2 Play-LMP Limitations

Several conceptual issues arise from the goal labeling strategy. Play-LMP labels the last state in the sampled play sequence as the goal for all earlier states (hindsight relabelling). The assumption embedded in this method is that the last state of any sequence of length H in play adequately represents the human's goal when choosing their actions (i.e., that the goal implies the actions). This can lead to the following issues relating to proper credit assignment:

Challenges of Fixed and Short Horizons: Inability to Capture Skills with Variable or Long Horizons: As discussed, choosing the horizon length is a balancing act between picking capturing longer horizon task goals and behaviors in the latent space (posterior) and ensuring the predictability of those skills from just the start and goal states at test time (prior). Due to this trade-off, Play-LMP can fail to capture the true goals when play consists of tasks whose horizon length is variable and/or large. In these settings, the longer horizon goal states will be out of distribution for the prior, and thus the policy can fail to produce the correct behavior. This trade-off in horizon lengths is visualized in Fig. 5 (a) for a pushing example.

Challenges of Short Horizons: Uninformative Goals: On the flip side, short sampled windows will often not contain any changes to the environment. One can imagine if the robot is re-orienting or servoing to an object, the goal state from the end of the short window will not be any different than the starting state, and thus uninformative for the prior network, even if the trajectory itself is

nontrivial – see Fig. 5 (b), where τ only captures the reaching phase of pushing. Thus Play-LMP might learn to map many sequences in play to the null prior, $E'(s,s^o)$, where s^o are the same environment state features in s.

Challenges of Fixed Horizons: Task Switching: Another issue related to sampling windows randomly from play is that we might sample a window that contains a logical boundary between two tasks. Thus, the labelled goal (last state in the window) belongs to the second task, and will not necessarily inform the behavior in the first task – see Fig. 5 (c), where τ contains both a pushing and lifting task back to back. Thus the prior will likely not be able to predict the correct plan distribution for the policy to use, and will incur a high KL penalty that might shape the latent space disproportionately.

Challenges of Random Trajectory Sampling: Even if we could perfectly label goals for each play trajectory τ , it might be the case that not all sub-sequences should be "equal" in terms of their contribution to the latent space structure. For example, a sequence involving a robot reaching an object does not contain as many critical states – states where the policy must be precise and accurate – as a sequence involving picking up a block or rotating it in-hand. With Play-LMP, both of these types of sequences would be weighted equally in the construction of the latent space. We posit that the latent space should attend more to critical states in play like object interactions; therefore, prioritizing encoding sequences from object interactions will enable a more useful and information rich latent plan space.

Fundamentally, the issues above relate to a failure of *credit assignment*, stemming from the short and fixed length of the sampled play sequences: the inferred goal does not actually represent the demonstrator's true goal. Our method PLATO considers tasks with variable horizon by leveraging object interactions with the environment. In the next section, we outline the implementation details for PLATO.

B PLATO Implementation Details

Next we will discuss the training procedure and low-level implementation details for PLATO, involving further algorithm details, network architectures, and hyperparameter choices.

B.1 Training Procedure

```
Algorithm 2 PLATO Training
1: Given: H^{(i)}, H^{(-)}, play data D_{\text{play}}, interaction criteria f^{(i)}, 2: D_{\text{play}}^{(-)}, D_{\text{play}}^{(i)}, D_{\text{play}}^{(+)} = f^{(i)}(D_{\text{play}})
                                                                                                                                                                                 ▷ Split into interactions
 3: Initialize E, E', \pi
 4: while not converged do
           \begin{split} \boldsymbol{\tau}^{(-)}, \boldsymbol{\tau}^{(i)}, \boldsymbol{\tau}^{(+)} &\sim \boldsymbol{D}_{\text{play}}^{(-)}, \boldsymbol{D}_{\text{play}}^{(i)}, \boldsymbol{D}_{\text{play}}^{(+)}, \\ \text{Sample } o_g &\sim \{o_t^{(+)}\} \end{split}
           p(z) \leftarrow E(\tau^{(i)})
                                                                                                                                                        ▶ Posterior Affordance Distribution
           p(z') \leftarrow E'(o_1^{(i)}, o_q)
 8:
                                                                                                                                                                 ▶ Prior Affordance Distribution
           z \sim p(z)
\tilde{a}_{1:H^{(i)}}^{(i)} \leftarrow \pi(s_{1:H^{(i)}}^{(i)}, o_g, z)
\tilde{a}_{1:H^{(-)}}^{(-)} \leftarrow \pi(s_{1:H^{(-)}}^{(-)}, o_g, z)
                                                                                                                                                                                   ⊳ Policy in Interaction
11:
                                                                                                                                                                           ⊳ Policy in Pre-Interaction
              Compute \mathcal{L}_{PLATO} with Eq. (1) and update \pi, E, E'.
```

We now outline the training procedure for PLATO in greater detail, with Algorithm 1 reproduced above in Algorithm 2 for convenience. First, we sample segmented pre-interaction, interaction, and post-interaction periods from the play dataset. We then sample fixed length windows $\tau^{(i)}=(s_1,a_1,...,s_{H^{(i)}})$ from the interaction period and $\tau^{(-)}=(s_1,a_1,...,s_{H^{(-)}})$ from the pre-interaction period (Line 5 in Alg. 2). Note that the post-interaction period is just the next pre-interaction period, and thus is still sampled for the next interaction. In the pull example in Fig. 1, $\tau^{(i)}$ is the pulling motion, and $\tau^{(-)}$ is the reaching motion before pulling. We sample fixed horizon snapshots of

each period for computational efficiency; however, the duration between $\tau^{(-)}$ and $\tau^{(i)}$ can vary tremendously, and so we are still able to capture variable and long horizon chains of events despite only sampling fixed horizon windows within each period.

As described in Sec. 3 in the main text, the goal object state o_g for this chain of events can be selected as any object state after the interaction period's last object state $o_{H^{(i)}}^{(i)}$ and before the end of the post-interaction period (Line 6 in Alg. 2). During the interaction to post-interaction range, we know that the object trajectory is determined only by the interaction window actions and obstacles in the scene. For example, if we grab a block and then launch it along the table, the post-interaction period will consist of the block sliding; any state along that slide directly results from grabbing and launching. Thus the goal o_g can correctly be attributed to affordance z.

Having sampled $\tau^{(-)}$, $\tau^{(i)}$, and o_g , PLATO learns a goal-conditioned policy to reproduce the actions in both $\tau^{(-)}$ and $\tau^{(i)}$. Our insight is that much of the diversity in task-relevant behavior is contained during the interaction period, so instead of encoding $\tau^{(-)}$ and $\tau^{(i)}$ separately, we only encode $\tau^{(i)}$ with posterior E (Line 7 in Alg. 2). Now, the policy during interaction is conditioned on the $z \sim E(\tau^{(i)})$. Importantly, the policy during pre-interaction also conditions only on z, representing the future interaction (Line 10-11 in Alg. 2).

B.2 Architecture

We follow a similar implementation as Play-LMP for the posterior, prior, and policy networks, as described in [3]. The posterior, prior, and policy networks are implemented as a Bidirectional GRU-RNN, an MLP, and a Unidirectional GRU-RNN, respectively. Input trajectories to the posterior include both robot and object state information, but aligning with Play-LMP we leave out actions for the posterior input, as including actions empirically worsens performance. Actions are target positions and orientations of the robot, since these are flexible and intuitive enough for humans to operate. All action reconstruction losses use Mean Absolute Error (deterministic actions), since empirically we found little benefit to using probabilistic actions with a negative log-likelihood loss. Specific architecture choices for each environment and method are detailed in Table 1, as determined by extensive hyperparameter sweeps. For PLATO, we set $H^{(i)} = H^{(-)} = H$, and pre-interaction reconstruction loss weight $\alpha = 1$. Included in this hyperparameter sweep are minor horizon variations as employed in Relay Policy Learning [4], which we found not to benefit policy learning for our settings.

PLATO additionally uses a "soft-boundary length" parameter (S) to allow for some flexibility in what is considered the boundary of interaction and pre-interaction during sampling. If c_s and c_e are the true contact start and end indices, then $\tau^{(i)}$ is sampled between c_s-S and c_e+S . Likewise, $\tau^{(-)}$ is sampled from between 0 and c_s+S . This creates overlap between the pre-interaction and interaction regions, which we found empirically is necessary such that the policy can be trained contiguously across the contact border. In practice, we set S=H/2, since this allows for full coverage of the contact border during sampling.

C Experimental Details

In this section we outline the environments, our real world setup, tasks, data collection, and evaluation process we employed. Each environment has substantial variability, and scripted policies are similarly designed to be quite diverse with sizeable injected noise.

C.1 Environments and Tasks

For all simulated environments, the contact signal used in simulation is the binary contact information between the robot and the rest of the scene, which can easily be computed in pybullet (3D) and pymunk (2D).

Block2D: The first environment is a 2D continuous block manipulation environment implemented with the PyMunk 2D physics engine [25]. Blocks in the environment are sized, massed, and positioned randomly at every reset. The ego agent (red) can interact with these blocks in the presence of gravity, with a special "tether" action that creates a link constraint if the ego agent is close enough to the block. This environment is meant as a 2D analog to more challenging block manipulations.

Environment	Method	β	H (seconds)	z	$\pi\text{-hidden}$	E-hidden	$E^\prime\text{-width}$
	Play-GCBC	N/A	2	N/A	128	N/A	N/A
	Play-LMP	1e-3	2	16	64	128	128
Block2D	PLATO	1e-3	2	16	64	128	128
DIOCK2D	Play-GCBC (H)	N/A	2	N/A	256	N/A	N/A
	Play-LMP (H)	1e-3	2	16	256	128	256
	PLATO (H)	1e-3	2	16	256	128	256
	Play-GCBC	N/A	4	N/A	128	N/A	N/A
3D-Flat	Play-LMP	1e-3	4	64	128	128	256
	PLATO	1e-4	4	64	128	128	256
	Play-GCBC	N/A	4	N/A	128	N/A	N/A
3D-Platforms	Play-LMP	1e-4	4	64	128	128	256
	PLATO	1e-4	3	64	128	128	256
	Play-GCBC	N/A	4	N/A	256	N/A	N/A
Mug-3D	Play-LMP	1e-4	4	64	256	256	256
	PLATO	1e-4	4	64	256	256	256
	Play-GCBC	N/A	4	N/A	256	N/A	N/A
Playroom3D	Play-LMP	1e-3	4	64	256	256	256
	PLATO	1e-4	4	64	256	256	256

Table 1: Best performing hyperparameters for each environment, method, and data source. β controls the regularization on the posterior from the learned prior. Each method is quite sensitive to this amount of regularization. H is the horizon length and controls the length of the sampled trajectory for the posterior encoder. All methods are quite sensitive to this as well. z is the latent vector, and |z| is the latent dimensionality. π -hidden and E-hidden are the hidden sizes for π and E respectively, and these control the expressiveness of each network. E'-width is the width of the prior network. We do not vary the number of layers in each network, which are chosen to be the same as in prior work.

As shown in Figure 6, we constructed a set of block manipulation primitives to evaluate on: Push, Pull, Lift, Tip, and Side-Rotate.

Block3D-Flat: The second environment is a 3D block manipulation environment, involving a simulated Franka Emika Panda robot arm, 3D blocks, and a table playground area surrounded by walls. Similarly to Block2D, the block dimensions, initial positions, and masses here are randomly sampled. The Panda robot arm uses an operational space torque controller to exert realistic and bounded forces on the blocks. In this environment, we implement two dimensional Push primitives along the table surface, as well as a Top-Rotate primitive to control the z-axis rotation of the block in either direction (see Figure 6). Results for this environment, shown in Table 4, were not included in the main text due to space limitations.

Block3D-Platforms: The third environment builds on Block3D-Flat, but introduces platforms along the walls to enable even more complex primitives like lifting and placing. Here, we implement two dimensional Push primitives, and Lift-Place primitives in all cardinal directions on the table plane (see Figure 6).

Mug3D-Platforms: The fourth environment is like the Block3D-Platforms environment, but uses a challenging mug object of varied size and mass, instead of blocks. Here, we implement two dimensional Push primitives, and Lift-Place primitives in all cardinal directions on the table plane, and additionally a Mug Rotate primitive similar to Block3D-Flat (see Figure 6 for an example of Mug Rotate). These primitives each require unique types of affordances for a mug, and this environment is designed to test how learning from play methods can adapt to more precise manipulation tasks.

Playroom3D: The fifth environment is similar to the one used in prior work, involving several dynamic objects: a block on the table, a drawer, a cabinet door, and buttons in the cabinet. This environment is challenging since it involves learning affordances over multiple objects, as well as

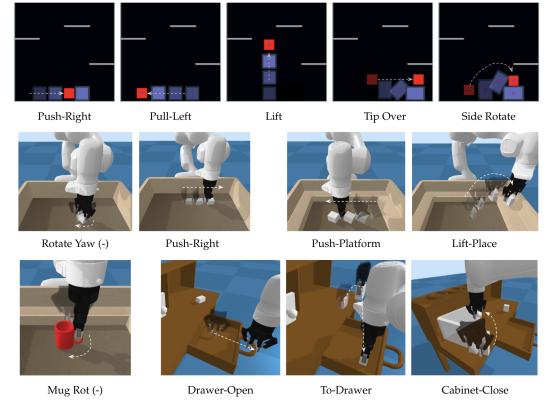


Figure 6: <u>First Row</u>: Block2D Environment Primitive Examples. Each of these primitives can be executed from a variety of object initial conditions, masses, and dimensions. <u>Second Row</u>: Block3D and Block3D-Platform Primitive Examples. Again, object initial conditions, masses, and dimensions are varied during play. The left two primitives shown are taken from **Block3D-Flat**, and the right two are taken from **Block3D-Platforms**. These represent a subset of the evaluation primitives in the 3D environment, and are meant to show the diversity of tasks and behaviors our method is evaluated on. <u>Third Row</u>: The left image shows an example primitive in **Mug3D-Platforms**. The right three images show sample tasks from **Playroom3D**.

interactions between them (like putting objects into the drawer or cabinet). For the cabinet, we test opening and closing actions, and likewise for the drawer. For the object, we test Pushing Left/Right primitives as well as moving the object To-drawer and To-cabinet (see Figure 6 for examples). These actions require many time-steps and have several bottleneck states. Play data also consists of moving the object From-drawer and From-cabinet, and button pressing (see Appendix D).

Real Robot Environment: Finally, we create a real robot environment that mirrors our simulation environment. Our setup is shown in Figure 7. We use three mounted Realsense SR300 cameras to robustly detect the pose of the green cube in the presence of occlusions from the robot, using OpenCV's Aruco tag detection for 3D pose estimation. As with our simulation experiments, we condition policies on full state of the environment (object 6D pose + end effector 6D pose + gripper state) rather than images. We add noise to the object state in simulation to match the noise seen in real world block state estimation. This along with an identical action space helps to reduce the sim-to-real gap and enables immediate deployment of simulation trained policies in the real world. Regardless, the dynamics of the object are still starkly different than those in our simulation dataset: the block is lightweight and slightly deformable, and has high friction with the table that can cause it to flip over instead of slide on the table. We test pushing primitives in this environment to demonstrate that PLATO is scalable to real world tasks with minimal data augmentation.

Real-World Deployment: To deploy this system more generically (including training on real world data), there are several key infrastructure additions beyond the setup we have shown here. In terms of hardware, the robot would need an additional contact sensing patch on the gripper or a force/torque sensor at the end effector to automatically detect interaction with the scene, as well as several cameras to observe the scene. Note that contact readings are only used during training (for the purpose

of interaction segmentation). Since our real world setup did not involve training on real world data, we did not need to add these additional sensors. However, we believe this modification should be quite straightforward. With pressure sensing (which is closest to what we do in simulation), there is a limitation that only interactions with the pressure sensing portion of the end effector will be counted during segmentation. In terms of software, a robust object 6D pose detection algorithm would be utilized to detect the object states under potential occlusion. With these additions, our method should readily scale to many real robot systems.

Example Task primitives for simulation evaluations are shown in Fig. 6. There is substantial withintask noise for scripted data to more closely resemble human data and real world conditions.

C.2 Task Complexity

We believe our set of tasks covers a wide spectrum of levels of complexity, from simpler pushing tasks to more complex door opening or object rotation tasks.

Firstly, we have incorporated diversity in primitives and variations in objects, which we emphasize is crucial and is not present in prior work such as Play-LMP. The tasks themselves cover a number of diverse primitives, where each "primitive" consists of variations in the goal, for example varying pushing distance or lifting placement location. The motions of the scripted primitives also have sizable variations in intermediate waypoints, speed, etc. Within all of our environments, we inject large variations in the positions, orientations, each size dimension, and masses of each of the blocks and mugs, which each require different strategies from the robot's perspective. Furthermore, grasping the mug involves a very precise interaction with the handle, which contrasts the wide, centered grasp used with blocks. This is in comparison to the closest prior work, i.e., Play-LMP tasks, which usually are projected to far fewer primitives and variations in the policy. The Play-LMP tasks largely involve either fixed object shapes with limited random pose initialization or static scene elements like buttons and constrained unchanging drawers. This might make it seem that these tasks are complex at the surface visual level, but we argue that our set of tasks and primitives require a much greater range of behaviors, such as grasping different shapes, rotating blocks to a wide spectrum of new orientations, and handling a variety of block masses, and thus these tasks are more complex. We would like to emphasize that visually interesting environments (e.g., added buttons or static objects as in Play-LMP), do not really add to the complexity of policy learning. What makes policy learning challenging is variations in behaviors and object properties, which we extensively test with our experiments. We believe that the performance of Play-LMP suffers in these settings precisely because the tasks are more complex for policy learning.

Secondly, our experiments include a visually interesting and complex environment, Playroom3D, which represents a more challenging version of the tasks used in Play-LMP involving some of the same underlying assets but having multiple randomized objects. To make the tasks even more challenging, we added the cabinet door opening and closing tasks. This "door opening" affordance was not included in the Play-LMP prior work, and is significantly harder to learn (see Play-LMP and Play-GCBC performance). We thus believe that our tasks are significantly more complex and diverse compared to prior work in this domain and adequately demonstrate the performance of our algorithm and a significant gap with prior work.

C.3 Evaluation Details

To evaluate each method, we evaluate across a set of environment-specific primitives by first running the primitive at the current state to generate a goal, then resetting and running the policy conditioned on that goal. We separately evaluate tasks by what we considered to be "similar" affordances. For example, push-left and push-right are considered as two separate tasks. However, there is sizeable variation in the set of goals that comprise each individual task. For example, with pushing, there is variation in the pushing distance, as well as all inherent variations in object properties in the environment (e.g., mass, shape, etc) discussed previously. Figure 8 provides some intuition that our affordance space learns to cluster by these directions. Success metrics are specified according to the task primitive being tested, usually involving the distance to either the goal object position or orientation. Evaluation times out after a fixed number of steps if the given method is unsuccessful. For stability, the latent vector z' is sampled once every $T \leq H$ steps and held constant during action decoding by the policy for the next T steps. For Play-LMP and PLATO, latent vectors are recomputed at the same frequency during evaluation, but the goal is fixed for the evaluation period.

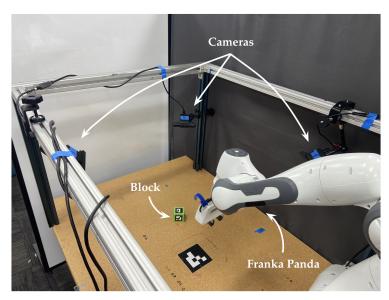


Figure 7: Block-Real environment. We use the Franka Emika Panda 7DOF robot arm for our experiments. The green cube we use for block manipulation tasks is shown in the middle, with ArUco tags on each face for pose detection. The 6D pose of the object is estimated via a multi-camera setup in order to be robust to occlusion by the robot, as shown here with the camera on the far right. The ArUco tag in the middle is used to calibrate the extrinsics of the cameras and localize the object frame of reference relative to the robot.

C.4 Data Collection

To evaluate our method across a wide variety of affordances available in the environment, we collect a large dataset of play data (roughly 15 hours) under artificial demonstrator agents that repetitively and randomly choose between the hand designed primitives for each environment. Proprioceptive state information includes robot end-effector 6D poses and twists, and object state information includes the object 6D poses, the 1D cabinet angle, and the 1D drawer open distance. In order to replicate human play as closely as possible, we add large amounts of variation to the primitives through randomized parameters. For example, with the pulling primitive, we vary the speed and duration of the pulling motion, as well as the reaching behavior. The boundaries between primitives are assumed to be unknown during training, like in human play data.

We also separately collect a smaller dataset (roughly 8 hours) of noisy human play data to evaluate the Block2D environments. From these experiments, we draw insights on the differences between human and scripted play data. Human play data for the Block2D task is collected using a keyboard control interface. Arrow keys control the ego agent position, while 'g' controls the grabbing tether action. One proficient user was used to collect the data, and was shown a set of tasks (the evaluation tasks) to perform during play with the instruction of trying to equally represent each task in their play, but they were also clearly informed that they were not limited to performing just these tasks. This user was given 15 minutes to practice in the environment, after which point data collection began. In future work, we hope to reduce the dependence on balanced, curated datasets and allow play to be truly freeform.

D Results and Analysis

In this section we discuss in greater depth our results, additional experiments, and tables with exact success rates to complement the bar-plot figures in the main text (Figures 2, 3, 4).

D.1 Detailed Results

Block2D: In Table 2, we see that PLATO substantially outperforms the baselines on each task. Play-LMP and Play-GCBC get roughly similar performance on most tasks, and struggle the most on tasks involving the tether action (Pulling and Side-Rotate). To further study the effects of encoding just the interaction period in the latent space, we implement PLATO-PRE, a version of PLATO

	Push-L	Push-R	Pull-L	Pull-R	Lift	Tip	Side-Rot
Play-GCBC	74.4(4.4)	84.5(2.9)	30.0(4.8)	18.6(3.4)	47.8(2.6)	72.8(3.8)	37.3(1.0)
Play-LMP	87.5(1.5)	89.9(1.0)	36.5(8.5)	17.9(4.0)	42.0(2.6)	81.0(2.0)	29.0(2.7)
PLATO-R	83.9(4.0)	86.6(6.12)	25.4(4.5)	38.1(0.5)	50.2(4.0)	79.4(2.6)	44.2(2.7)
PLATO-PRE	95.9(1.9)	95.7(1.0)	66(11.9)	67.1(7.4)	80.9(7.1)	84.3(1.6)	59(0.9)
PLATO	99.1(0.5)	98.8(1.2)	69.0(3.8)	81.4(3.0)	71.5(2.3)	86.9(3.6)	73.8 (1.0)
Play-GCBC (H)	33.3(5.1)	52.4(3.0)	21.1(11.1)	49.3(3.4)	40.7(8.6)	52.0(4.7)	34.1(3.3)
Play-LMP (H)	54.8(4.7)	62.2(4.4)	27.9(7.4)	58.6(1.6)	29.2(7.2)	53.9(3.7)	29.2(1.2)
PLATO (H)	81.3(3.0)	82.2(1.7)	65.3(2.7)	69.5(0.5)	65.6(2.8)	79.9 (0.6)	54.0 (1.7)

Table 2: Block2D Success Rates in percentages in the form mean(std-err), trained over 3 random seeds and evaluated on various Push, Pull, Lift, Tip, and Side-Rotate primitives. Block sizes are randomized in each dimension, and blocks are initialized in random positions along the bottom of the grid. Our method PLATO outperforms all prior methods on both scripted and human data. PLATO-PRE includes pre-interaction information in the learned latent space (amounting to passing both $\tau^{(i)}$ and $\tau^{(-)}$ into the posterior E), but increases the training time compared to PLATO. PLATO-R incorporates the current robot state into the prior, representing the non object-centric version of PLATO. Object-centric methods that learn from interactions (PLATO, PLATO-PRE) perform much better than their counterparts.

	Push-L	Push-R	Pull-L	Pull-R	Lift	Tip	Side-Rot
PLATO	99.1(0.5)	98.8(1.2)	69.0(3.8)	81.4(3.0)	71.5(2.3)	86.9(3.6)	73.8(1.0)
PLATO-FC(4%)	100	95.9	61.3	57.0	78.4	93.0	68.3
PLATO-FC(8%)	97.0	96.7	35.0	65.6	80.4	95.6	48.7
PLATO-FC(12%)	94.8	100	61.5	10.1	60.8	92.3	63.6

Table 3: Contact Signal Ablation Experiment. Here we artificially add fake contact signals outside of interactions (e.g., during pre-interaction), causing false interactions to be segmented during training. We evaluate PLATO-FC(%), where % denotes the percentage of *interactions* that are false positives. We show results for a single seed of each PLATO for 4%, 8%, and 12%. Pulling tasks have some variance in performance under added false contact, however considering how much data is affected, PLATO is quite robust for all tasks.

in which the posterior encodes sampled trajectories from both the pre-interaction period $\tau^{(-)}$ and interaction period $\tau^{(i)}$, instead of just the interaction period. We see that PLATO and PLATO-PRE do roughly equivalently on average, with PLATO doing substantially better on Side-Rotate, but PLATO-PRE doing better on the lift primitive. Note that PLATO-PRE is slower to train since the posterior recurrent encoder receives a much longer sequence. We hypothesize that this is due to the tradeoff of various tasks between the complexity of the interaction and the complexity of pre-interaction behaviors. Overall, we can conclude that the including pre-interaction trajectories in the latent space (PLATO-PRE) is not uniformly better than only including interaction trajectories (PLATO), and thus does not warrant the added training time. This confirms our intuition that for complex interaction sequences like Side-Rotate, the interaction period contains enough information from the perspective of representation learning. Importantly, we see that framing play data through the lens of object interactions (PLATO, PLATO-PRE) results in much better policy learning than prior state-of-the-art methods.

Interestingly, performance for all methods is worse on the human generated data than scripted data. We attribute this to the fact that despite the significant noise added to the scripted primitives during data collection, scripted data still has cleaner and more successful demonstrations of each task than human data, which can contain many sub-optimal trajectories due to the challenges of teleoperation. For example, we observed that the Side-Rotate primitive is only successful in the human play dataset around 70% of the time. Additionally, humans can demonstrate the same task in many more ways than we could possibly script (e.g., by elongating the duration of a pulling motion or picking a wildly different spot to pull from), resulting in much more complex plans and policies to learn. This is supported by the fact that, as shown in Table 1, the best performing architectures for each method on human play data involved larger policy hidden sizes as compared to the best performing methods on scripted play data.

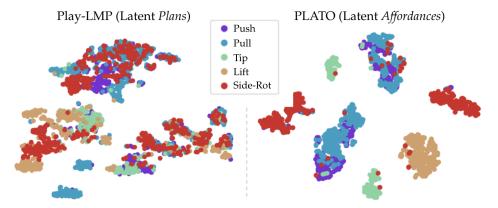


Figure 8: Learned latent spaces for Play-LMP (left) and our method PLATO (right) in the Block2D environment, plotted in 2 dimensions (from original 16 dimensions) with t-SNE. While there is overlap between tasks in the learned plan space in Play-LMP, tasks are much more separable with the learned affordance space in PLATO. We see separate clusters for each primitive likely relating to the direction of the primitive, and we also see within-cluster variation for each primitive. Note that push and pull have sizeable overlap since these represent similar motions (motion left and right on the ground). This clustering suggests that learning from interactions helps the posterior E to provide simpler and more informative latent representations for the policy.

For these 2D environments, we additionally examine the structure of the latent space learned by PLATO compared to prior work. Figure 8 shows example learned latent spaces for both Play-LMP and our method PLATO on Block2D; we see that by training on interactions, PLATO recovers a latent space that is more separable by task than Play-LMP despite not being presented with any task labels during training.

Additionally, in Table 3), we ablate the quality of the interaction signal (contact) in order to determine the important of clean contact signals. To accomplish this, we add in *false positive* contact signals to the data (e.g., during pre-interaction). This causes some percentage of the interactions sampled during training (denoted by FC(%) in Table 3) to actually be non-contact sequences. We see that PLATO is quite robust to added false positive contacts. The Pull task has high variance in performance here under added contacts, however there is no clear relationship between more contact noise and worse performance for any of the tasks. This is likely due to the fact that all models condition on the start and goal states, and false contact signals usually involve no change in the object state; therefore, from the model's perspective, sampled false interactions will be completely distinguishable from the task affordances that we care about at test time.

					Rotate (+)	` '
GCBC	71.3(3.0)	65.7(10.0)	57.5(11.2)	50.1(7.8)	53.7(2.6)	42.4(1.9)
LMP	68.0(4.5)	63.0(4.4)	50.6(2.3)	28.0(9.0)	39.6(3.2)	40.7(3.4)
PLATO	77.4(4.3)	89.4(5.5)	74.0(4.7)	84.0(5.0)	83.3(4.4)	78.6 (6.7)

Table 4: Block3D-Flat Success Rates in percentages in the form mean(std-err), trained over 3 random seeds and evaluated on pushing and rotation tasks. Again, our method PLATO outperforms all prior methods, especially on the harder rotation primitives. Interestingly, we also see that Play-LMP does not do as well as Play-GCBC on several primitives, likely due to the prior and policy being out of distribution at test time.

Block3D-Flat: Results for this block pushing and rotating environment (not presented in the main text due to space limitations) are shown in Table 4. Here again, PLATO is able to significantly outperform the baselines on all of the tasks, especially in the more fine-grained rotation motions, as well as on the pushing tasks.

Block3D-Platforms: Table 5 shows the results for Block3D-Platforms. PLATO outperforms the baseline methods in the pushing tasks, but the difference is substantially greater for the lifting tasks. These lifting tasks have longer, variable horizons and involve several bottleneck states (e.g., securely

			Push-L						Lift-R
Ī			69.7(10.4)						
	LMP	84.7(7.3)	85.3(3.4)	90(2.9)	91.2(2.6)	47.2(5.6)	55.1(4.6)	48.6(10.9)	51.8(13.6)
	PLATO	98.3(0.4)	97.2(1.4)	97.8(2.2)	99.7(0.3)	75.6 (1.9)	94.0(2.0)	81.0(1.1)	92.8(1.7)

Table 5: Block3D-Platform Success Rates in percentages in the form mean(std-err), trained over 3 random seeds and evaluated on Lift-Place and Push primitives. PLATO is the only method able to do well on all evaluation primitives and do so consistently, for a wide variety of object dimensions and initial conditions.

grasping, clearing the platform, dropping on the platform). PLATO is the only method capable of performing well on both the longer horizon tasks (lifting) and shorter horizon tasks (pushing).

	Push-F	Push-L	Push-B	Push-R	Rot(+)	Rot(-)	Lift-F	Lift-L	Lift-B	Lift-R
GCBC	83.7(1.3)	91.3(1.5)	86.3(5.4)	93.0(2.5)	78.7(5.8)	81.0(0.6)	11.1(6.4)	27.7(2.7)	4.4(2.5)	20.6(1.3)
LMP	85.0(2.6)	87.6(3.8)	89.9(0.6)	90.6(3.1)	51.3(5.2)	63.3(3.2)	59.3(5.8)	74.3(4.7)	51.7(6.6)	64.8(5.3)
PLATO	94.7(0.9)	92.0(2.1)	97.2(1.0)	98.7 (0.7)	81.3(2.5)	86.3(6.1)	81.7(0.3)	84.3(3.0)	76.7(10)	85.2(3.6)
LMP(S)	100	100	92.7	100	84.8	85.5	67.6	70.9	70.1	87.7
PLATO(S)	100	100	97.4	100	100	100	92.3	90.2	82.4	85.6
LMP(S+)	66.8	76.0	75.8	64.6	62.5	87.1	47.3	44.8	35.3	34.9
PLATO(S+)	88.3	87.1	77.4	90.9	87.4	86.7	64.4	71.0	52.4	72.3

Table 6: Mug3D-Platform Success Rates in percentages in the form mean(std-err), trained over 3 random seeds and evaluated on Lift-Place, Rotate, Push primitives. In the first block of the table, we see that PLATO is the only method able to do well on all evaluation primitives and do so consistently, for a wide variety of object dimensions and initial conditions. In the second block of this table (generalization experiments, one seed), (S) denotes that the method was trained on a subset of initial mug orientations (simpler tasks). While in principle play will reach other mug orientations, this still skews the distribution of mug orientations towards the initial predefined set. (S+) denotes methods trained on these subset of initial orientations (same models as S) but evaluated on the full range of mug orientations as used in the first block. We see that while LMP performs close to PLATO within distribution (subset of initial mug orientations, significantly less task diversity), PLATO is much more robust than LMP when presented with the full initial object state distribution at test time, showing the generalization capacity of PLATO.

Mug3D-Platforms: Table 6 contains both the results for Mug3D-Platforms presented in the main text, as well as an additional ablation experiment testing the generalization capacity of PLATO in this environment. The main experiments (first block of Table 6) show that similar to in Block3D-Platforms, PLATO outperforms the methods on all tasks, where the gap is less stark for pushing tasks but especially large on the lifting tasks. For the rotation tasks, interestingly Play-GCBC outperforms Play-LMP. We speculate that in cases like this, the Play-LMP policy might be too dependent on the plan z, and thus suffers at test time when the prior outputs an approximate z distribution given only partial information. In contrast, Play-GCBC is substantially worse than Play-LMP for the lifting task. This suggests that task variability is much larger for lifting than rotating and pushing, and hence the latent plan in Play-LMP helps the policy disambiguate between this variability. Overall, PLATO performs better and more consistently on all the tasks than either Play-GCBC or Play-LMP.

The second set of experiments (second block in Table 6) illustrate the robustness of PLATO to state/action/goal distribution shift. We train Play-LMP and PLATO on a subset of initial mug orientations in the mug environment. Specifically, the initial randomized z-axis orientation (yaw) of the mug at the start of each episode will be just a 90 degree cut of the full 360 degree range. While in principle sequential play will be able to eventually see tasks demonstrated for orientations outside this range, this initialization greatly skews the distribution of mug orientations for all demonstrated tasks towards the starting set of orientations. The first two rows of the second block in Table 6 (S) show the performance of Play-LMP and PLATO when evaluated on the tasks used for training (90 degree cut for initial mug orientations). We see here that due to the lower variability in tasks, Play-LMP performs quite well within distribution, notably on the lifting tasks (in contrast to the results from the first block in the table). PLATO performs better than Play-LMP across all tasks, consistent with the results in the first block, although the gap is reduced due to limited task variability. The second two rows of the second block in Table 6 (S+) show the performance when the same models from the first two rows (S) are evaluated on the full swath of initial mug orientations.

	Push-L	Push-R	Cab-O	Cab-C	Dr-O	Dr-C	To-Cab	To-Dr	Fr-Cab	Fr-Dr	Btn1	Btn2
GCBC	29.1(3.5)	52.9(4.0)	47.7(21.4)	1.2(1.2)	86.7(6.4)	22.5(11)	52.2(12)	60.6(10)	6.3(4.1)	3.7(0.3)	68.7(11)	66.7(18)
LMP	25.5(2.5)	60.3(10.5)	61.2(8.3)	0(0)	65.3(13)	27.2(13)	34.7(6.2)	45.5(7.1)	1.0(1.0)	19.7(6.5)	63.2(9.2)	42.3(13)
PLATO	76.3(15)	66.7(15)	100(0)	78.7(6/9)	100(0)	100(0)	77.3(3.7)	60.4(2.0)	11.7(1.8)	58.3(3.5)	70.9(2.1)	100(0)

Table 7: Playroom3D Success Rates in percentages in the form mean(std-err), trained over 3 random seeds and evaluated on Push, Cabinet, Drawer, To/From-Cabinet/Drawer, and Button Pressing primitives. PLATO is the only method able scale to all evaluation primitives, for a wide variety of object sizes and initial conditions.

As explained previously, both Play-LMP and PLATO have seen a limited set of tasks with these orientations due to the sequential nature of play, but only PLATO is able to retain good performance across all tasks. This demonstrates that by biasing the latent space towards learning object affordances, PLATO is better able to capture the full distribution of demonstrated behaviors, even those infrequently demonstrated, and thus is more adept under distribution shift in the test time tasks.

Playroom3D: In Table 7, we show the results on Playroom3D. Surprisingly, pushing tasks are much harder in this environment, which we speculate is due to the presence of multiple objects and thus a higher variability in how objects can be interacted with. PLATO is able to retain good performance on the pushing tasks, in contrast to the baselines. For opening and closing the cabinet and drawer, we see that PLATO is able to do quite well, but Play-LMP and Play-GCBC perform quite poorly. There is quite a diversity in horizon lengths across these different tasks and different instantiations of each task, which we believe contributes to this large gap. For the Cabinet closing task, there are several critical bottleneck states, for example being able to servo the arm above and to the other side of the cabinet door to reach the cabinet handle, as well as precisely grasping the handle. PLATO is the only method that learns to consistently perform this task. Likewise for the drawer tasks, PLATO gets 100% success for all random seeds, while performance is substantially worse on the baselines. We also evaluated two additional challenging tasks in the Playroom3D environment that were demonstrated during play less frequently: From-Cabinet and From-Drawer (pull object out of open cabinet, lift it out of open drawer), also with substantial object and primitive variation. The performance on these tasks are lower compared to the other tasks as they require many preconditions and thus are not equally represented during our collected play data. From-Cabinet also requires a novel end effector orientation and grasping procedure in order to avoid collision with the cabinet and table walls. We speculate that due to the novel motion, limited examples, and the presence of other tasks in the data, all methods do notably worse on the From-Cabinet task, however From-Drawer performance for Play-LMP and PLATO is closer to To-Drawer performance since these tasks involve similar object lifting behaviors. However, even with fewer examples in a crowded dataset of other tasks, there is still a large gap between PLATO and the next best method. Additionally, we evaluate on two button pressing tasks, where the buttons are in the cabinet space and thus reaching involves avoiding the cabinet door. Here, the period of contact is relatively short, but PLATO is able to reliably achieve higher success in these tasks as well. Overall, this environment demonstrates that PLATO gracefully scales to more complex environments with more diversity in object affordances and robot behaviors.

BlockReal: Results for our real world experiments are shown in Table 8. Both Play-LMP and PLATO are trained entirely in simulation. We add minor data augmentation in simulation in the form of gaussian noise for the object state estimates, to match the observed noise using our real world multi-camera object state estimation infrastructure. Note that both models get 90%+ success in simulation on each task, since these tasks have relatively low behavior and object diversity. We see that when deployed on the real world setup with no additional data, PLATO experiences only a minor performance degradation, suggesting that learning a latent *affordance* space is more robust than learning a latent *plan* space. These results are consistent with what we see in the Mug3D-Platforms generalization experiments in Table 6, however here we are testing generalization to entirely unseen real world physics, in contrast to task distribution shift in those experiments.

D.2 Additional Analysis

Variation in Performance for Similar Tasks: For several environments, semantically similar tasks like push-left and push-forward seem to have notably different success rates across many methods. Interestingly, those differences are mainly across different "axes" of the task, for example push-left and push-right usually have similar performance, and push-forward and push-backward also have

	Push-Left	Push-Back	Push-Right	Push-Forward
LMP	8/10	6/10	8/10	7/10
PLATO	8/10	8/10	9/10	10/10

Table 8: Results for BlockReal on pushing tasks. These models are trained entirely in simulation with minor data augmentation. We evaluate these models on a real robot setup, and see that performance degrades less for PLATO than for Play-LMP when presented with the real world object and robot dynamics. Methods that use play data are robust to environment changes, consistent with results from prior work [3].

similar performance, but these two sets have a gap in performance. We speculate that this is because of biases present in the data or the model that favor one axis over another, for example the relative presence of each task, or environmental difficulties in performing that task.

Quality of Interaction Segmentation: In all our experiments, we are not assuming access to perfect interaction segmentation. In fact, all of our environments will sometimes have imperfect interaction signals due to the demonstrator having notable noise (even in scripted policies, which have added noise) – for example, brushing against the table, object, or cabinet door on the way to perform a different task. Intermittent contact is also quite common in all of our play data. However, the smoothing on top of the interaction signal tends to clean many of these signals. See the Contact Ablation Experiments in Appendix D.1 for more analysis.

Furthermore, we pose this question: what defines "perfect" segmentation? In the framing of our method, any interaction with the environment, even accidental ones, are still valid for the affordance space to learn. If we accidentally brush the top of the cabinet door on our way to push an object, and the door opens slightly, this can be seen as a successful cabinet slight-open task. Since the start and goal object state are unique for this task, in theory it should not at all affect the affordance learning for a different start and goal object state. However, accidental interactions will start to affect learning if these interactions are common and bias the policy towards unsafe regions of the state (for example, if repeatedly brushing the top of the cabinet door on the way to push the block biases the policy away from pushing the block properly).

E Additional Limitations

We will now present a longer discussion of limitations presented in Sec. 5, as well as some additional limitations, to help guide future work.

Multi-Object Scenarios: As noted in Sec. 5, multi-object interaction scenarios like tool-use represent a key challenge for future work. While we show PLATO operating in multi-object environments in this work, we do not extend PLATO to multi-object interactions involving dynamic objects, for example hitting a puck with a hockey stick. In these settings, it might be difficult to obtain signals for interaction between the dynamic objects. We will give a couple of ideas of how future work might tackle this problem in the hopes of opening up a broader discussion. Before these ideas, one general point of clarification: PLATO introduces detecting "interaction" as a more general concept, which applies even when no contact occurs between the robot and the desired object (e.g., tool use). We used contact as our interaction signal since it was the most readily available for singleobject interactions. However for multi-object interactions, the notion of interaction still exists and our method still applies if we can detect these interactions somehow. Since PLATO only requires detecting interaction during training, one option is to have a human label interaction segments in their training data manually. If this is too time intensive, we can potentially leverage learned binary signals for interaction using limited supervised interaction labels. Consider the tool-use task of hitting a hockey puck with a hockey stick. Even though we cannot directly observe a contact signal between the stick and the puck, future work could learn to predict the interaction signal from visual (e.g. observing the stick hit the puck) and maybe even haptic information (force feedback of the hockey stick on the end-effector when hitting), using supervised labels. In addition, recent approaches like ComPILE [26] have learned to segment skills without any labels, and could be used to isolate an interaction signal for dynamic multi-object interactions in a self-supervised fashion. While designing such a system for detecting interaction might require some effort, we believe that PLATO's results suggest that such effort can result in serious performance gains for policy learning from play.

With this general notion of interaction in mind, we would like to present two potential ideas for future work to expand PLATO to multi-object scenarios:

- 1. One idea would be to learn an *embodiment-specific* affordance space. Then, each tool can be viewed as a different embodiment of the robot, and with knowledge of what tool is currently being used, we can learn affordances specific to that tool. At test time if we know what tool we picked up, the policy can leverage the affordance space of this particular tool in a manner similar to PLATO.
- 2. As another related idea, we could attempt to learn multiple *degrees* of interaction: e.g. similar to how PLATO learns the relationship between an object affordance (hockey stick moves) and robot skill (grasp and swing a stick), we might also learn the relationship between a second order object affordance (hockey puck moves) and a first order object affordance (hockey stick strikes). Then at test time the policy could reason backwards in time, first inferring the correct second order affordance (hockey puck moves), then the first order affordance (hockey stick swings), then the robot action to take (grasp and swing the stick).

We see our work as a first step towards exploring some of these interesting paradigms for interaction and multi-level skill reasoning.

Regularization Weight: As one might expect, both Play-LMP and our algorithm PLATO share many of the challenges of variational auto-encoders [27]. Recent replications of this work have shown that the regularization weight β has a sizeable effect on the final policy reconstruction error [28]. High values of β can yield posterior collapse of the latent space, where the plan posterior outputs distributions for differing trajectories that do not reflect these differences; low values of β can yield low reconstruction losses, but conversely cause the posterior plans to encode information that is hard to predict from just the start and end states. As a result, the learned prior may not match plans encoded by the posterior, thus hurting the policy. Future work might look into more expressive learned priors, such as mixture models, in order to better match the posterior and thereby reduce the sensitivity to β . Another direction could be finding alternate ways to specify tasks at test time, for example giving some notion of *how* a goal should be reached.

Human Sub-Optimality: Additionally, when collecting play data, certain challenging primitives attempted by humans might fail often. We noticed that humans often fail at Side-Rotate in Block2D, for example, and the resulting demonstration might look like a sub-optimal Push from the perspective of the posterior and the prior. This introduces even more plan variability into the latent space for the Push task, and thus hurts test time performance. Another interesting direction would be to better understand how sub-optimality affects the learned latent space and potentially develop a notion of trajectory "quality" to bias this latent space.

Use of Object State Estimation: As mentioned in Sec. 5, PLATO makes use of object state estimates when learning object-centric affordances. A natural question is how this work can be adapted to work with pure image inputs. Critically, the only obstacle to extending our method to image inputs is our learned prior network, which utilizes just ground-truth object state information rather than the full proprioceptive state and object state. Note that our PLATO-R ablation can be extended to learn from images without any additional modifications, however this ablated method lacks the benefits of object-focused affordance learning (see Section 4). While we did not explore learning from images, we believe our method can readily scale to images with a few modifications. One method would be to mask out only the object(s) of interest from the start and goal images before passing them into the prior to encourage a robot agnostic latent space. Another method would be to learn an object representation directly from images that is independent of the robot state (e.g. with contrastive learning on negative examples of different robot poses, but identical environment states).

Regardless, we showed that learning object centric representations actually improves the robustness of policies at test time, and thus justifies the extra effort of designing these representations. Furthermore, there is a sizable research field devoted to improving object state estimation using learning and filtering techniques (even involving state estimation under clutter), so we believe that object state estimation methods will get even more practical in the coming years.