

# Zero-shot Policy Learning with Spatial Temporal Reward Decomposition on Contingency-aware Observation

Huazhe Xu<sup>1\*</sup>, Boyuan Chen<sup>1\*</sup>, Yang Gao<sup>2</sup>, Trevor Darrell<sup>1</sup>

**Abstract**—It is a long-standing challenge to enable an intelligent agent to learn in one environment and generalize to an unseen environment without further data collection and finetuning. In this paper, we consider a zero shot generalization problem setup that complies with biological intelligent agents’ learning and generalization processes. The agent is first presented with previous experiences in the training environment, along with task description in the form of trajectory-level sparse rewards. Later when it is placed in the new testing environment, it is asked to perform the task without any interaction with the testing environment. We find this setting natural for biological creatures and at the same time, challenging for previous methods. Behavior cloning, state-of-art RL along with other zero-shot learning methods perform poorly on this benchmark. Given a set of experiences in the training environment, our method learns a neural function that decomposes the sparse reward into particular regions in a contingency-aware observation as a per step reward. Based on such decomposed rewards, we further learn a dynamics model and use Model Predictive Control (MPC) to obtain a policy. Since the rewards are decomposed to finer-granularity observations, they are naturally generalizable to new environments that are composed of similar basic elements. We demonstrate our method on a wide range of environments, including a classic video game – Super Mario Bros, as well as a robotic continuous control task. Please refer to the project page for more visualized results.<sup>1</sup>

## I. INTRODUCTION

While deep Reinforcement Learning (RL) methods have shown impressive performance on video games [1] and robotics tasks [2], [3], they solve each problem *tabula rasa*. Hence, it will be hard for them to generalize to new tasks without re-training even due to small changes. However, humans can quickly adapt their skills to a new task that requires similar priors *e.g.* physics, semantics and affordances to past experience. The priors can be learned from a spectrum of examples ranging from perfect demonstrative ones that accomplish certain tasks to aimless exploration.

A parameterized intelligent agent “Mario” who learns to reach the destination in the upper level in Figure 1 would fail to do the same in the lower new level because of the change of configurations and background, *e.g.* different placement of blocks, new monsters. When an inexperienced human player is controlling the Mario to move it to the right in the upper level, it might take many trials for him/her to realize the falling to a pit and approaching the “koopas”(turtle) from the left are harmful while standing on the top of the “koopas”(turtle) is not. However, once learned, one can

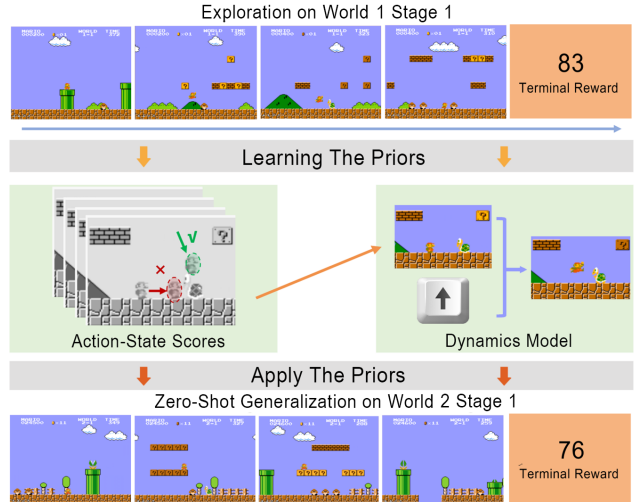


Fig. 1: Illustrative figure: An agent is learning priors from exploration data from World 1 Stage 1 in Nintendo Super Mario Bros game. In this paper, the agent focuses on learning two types of priors: learning an action-state preference score for contingency-aware observation and a dynamics model. The action-state scores on the middle left learns that approaching the “Koopas” from the left is undesirable while from the top is desirable. On the middle right, a dynamics model can be learned to predict a future state based on the current state and action. The agent can apply the priors to a new task World 2 Stage 1 to achieve reasonable policy with zero shot.

infer similar mechanisms in the lower level in Figure 1 without additional trials because human have a variety of priors including the concept of object, similarity, semantics, affordance, etc [4], [5]. In this paper, we teach machine agents to realize and utilize useful priors from exploration data in the form of decomposed rewards to generalize to new tasks without fine-tuning.

To achieve such zero-shot generalizable policy, a learning agent should have the ability to understand finer-granularity of the observation space *e.g.* to understand the value of a “koopas” in various configuration and background. However, these quintessentially human abilities are particularly hard for learning agents because of the lack of temporally and spatially fine-grained supervision signal and contemporary deep learning architectures are not designed for compositional properties of scenes. Many recent works rely heavily on the generalization ability of neural networks learning algorithms without capturing the compositional nature of scenes.

In this work, we propose a method, which leverages imperfect exploration data that only have terminal sparse rewards to learn decomposed rewards on specific regions of

<sup>1</sup> UC Berkeley

<sup>2</sup> Tsinghua University

\* Equal Contribution

<sup>1</sup><https://sites.google.com/view/sapnew/home>

an observation and further enable zero-shot generalization with a model predictive control (MPC) method. Specifically, given a batch of trajectories with terminal sparse rewards, we use a neural network to assign a reward for the contingency-aware observation  $o_t$  at timestep  $t$  so that the aggregation of the reward from each contingency-aware observation  $o_t$  can be an equivalence of the original sparse reward. We adopt the contingency-aware observations [6] that enables an agent to be aware of its own location. Further, we divide the contingency-aware observation into  $K$  sub-regions to obtain more compositional information. To further enable actionable agents to utilize the decomposed score, a neural dynamics model can be learned using self-supervision. We show that how an agent can take advantage of the decomposed rewards and the learned dynamics model with planning algorithms [7]. Our method is called SAP where ‘‘S’’ refers to scoring network used to decompose rewards, ‘‘A’’ refers to the aggregation of per step rewards for fitting the terminal rewards, and ‘‘P’’ refers to the planning part.

The proposed scoring function, beyond being a reward function for planning, can also be treated as an indicator of the existence of objects that affect the evaluation of a trajectory. We empirically evaluate the decomposed rewards for objects extracted in the context of human priors and hence find the potential of using our method as an unsupervised method for object discovery.

In this paper, we have two major contributions. First, we demonstrate the importance of decomposing sparse rewards into temporally and spatially smaller observation for obtaining zero-shot generalizable policy. Second, we develop a novel instance that uses our learning-based decomposition function and neural dynamics model that have strong performance on various challenging tasks.

## II. RELATED WORK

**Zero-Shot Generalization and Priors** To generalize in a new environment in a zero-shot manner, the agent needs to learn priors from its previous experiences, including priors on physics, semantics and affordances. Recently, researchers have shown the importance of priors in playing video games [5]. More works have also been done to utilize visual priors in many other domains such as robotics for generalization, etc. [8], [9], [10], [11], [12], [13], [14], [15] explicitly extended RL to handle object level learning. While our method does not explicitly model objects, we have shown that meaningful scores are learned for objects enabling SAP to generalize to new tasks in zero-shot manner. Recent works [16], [17] try to learn compositional skills for zero-shot transfer, which is complementary to the proposed method.

**Inverse Reinforcement Learning.** The seminal work [18] proposed inverse reinforcement learning (IRL). IRL aims to learn a reward function from a set of expert demonstrations. IRL and SAP fundamentally study different problem — IRL learns a reward function from *expert* demonstrations, while our method learns from exploratory data that is not necessarily related to any tasks. There are some works

dealing with violation of the assumptions of IRL, such as inaccurate perception of the state [19], [20], [21], [22], or incomplete dynamics model [23], [19], [24], [25], [26]; however, IRL does not study the case when the dynamics model is purely learned and the demonstrations are suboptimal. Recent work [27] proposed to leverage failed demonstrations with model-free IRL to perform grasping tasks; though sharing some intuition, our work is different because of the model-based nature.

**RL with Sparse Reward** When only sparse rewards are provided, an RL agent suffers a harder exploration problem. Previous work [28] studied the problem of reward shaping, *i.e.* how to change the form of the reward without affecting the optimal policy. The scoring-aggregating part can be thought as a novel form of learning-based reward shaping. A corpus of literature [29], [30], [31] try to learn the reward shaping automatically. However, the methods do not apply to the high-dimensional input such as image. One recent work RUDDER [32] utilizes an LSTM to decompose rewards into per-step rewards. This method can be thought of an scoring function of the full state in our framework.

There are more categories of methods to deal with this problem: (1) Unsupervised exploration strategies, such as curiosity-driven exploration [33], [34], or count-based exploration [35], [36] (2) In goal-conditioned tasks one can use Hindsight Experience Replay [37] to learn from experiences with different goals. (3) defining auxiliary tasks to learn a meaningful intermediate representations [38]. In contrast to previous methods, we effectively convert the single terminal reward to a set of rich intermediate representations, on top of which we can apply planning algorithms. Although model-based RL has been extensively studied, none of the previous work has explored the use of reward decomposition for zero-shot transfer.

## III. PROBLEM STATEMENT AND METHOD

### A. Problem Formulation

To be able to generalize better in an unseen environment, an intelligent agent should understand the consequence of its behavior both spatially and temporally. In the RL terminology, we propose to learn rewards that correspond to observations spatially and temporally from its past experiences. To facilitate such goals, we formulate the problem as follows.

Two environments  $\mathcal{E}_1, \mathcal{E}_2$  are sampled from the same task distribution.  $\mathcal{E}_1$  and  $\mathcal{E}_2$  share the same physics and goals but different configurations. e.g. different placement of objects, different terrains.

The agent is first presented with a bank of exploratory trajectories  $\{\tau_i\}, i = 1, 2 \dots N$  collected in training environment  $\mathcal{E}_1 = (\mathcal{S}, \mathcal{A}, p)$ . Each  $\tau_i$  is a trajectory and  $\tau_i = \{(s_t, \mathbf{a}_t)\}, t = 1, 2 \dots K_i$ . These trajectories are random explorations in the environment. Note that the agent only learns from the bank of trajectories without further environment interactions, which mimics human utilizing only prior experiences to perform a new task. We provide a scalar

terminal evaluation  $r(\tau)$  of the entire trajectory when a task  $\mathcal{T}$  is specified.

At test time, we evaluate task  $\mathcal{T}$  using zero extra interaction in the new environment,  $\mathcal{E}_2 = (\mathcal{S}', \mathcal{A}, p')$ . We assume identical action space. There is no reward provided at test time. In this paper, we focus on locomotion tasks with object interactions, such as Super Mario running with other objects in presence, or a Reacher robot acting with obstacles around.

### B. Spatial Temporal Reward Decomposition

In this section, we introduce the method to decompose the terminal sparse reward into specific time step and spatial location. First, we introduce the temporal decomposition and then discuss the spatial decomposition.

**Temporal Reward Decomposition** The temporal reward decomposition can be described as  $S_\theta(W(s_t), \mathbf{a}_t)$ . Here  $\theta$  denotes parameters in a neural network,  $W$  is a function that extracts contingency-aware observations from states. Here, contingency-aware means a subset of spatial global observation around the agent, such as pixels surrounding a game character or voxels around end-effector of a manipulator. We note that the neural network’s parameters are shared for every contingency-aware observation. Intuitively, this function measures how well an action  $\mathbf{a}_t$  performs on this particular state, and we refer to  $S_\theta$  as the *scoring function*.

To train this network, we aggregate the decomposed rewards  $S_\theta(W(s_t), \mathbf{a}_t)$  for each step into a single aggregated reward  $J$ , by an aggregating function  $G$ :

$$J_\theta(\tau) = G_{(s_t, \mathbf{a}_t) \in \tau} (S_\theta(W(s_t), \mathbf{a}_t))$$

The aggregated reward  $J$  are then fitted to the sparse terminal reward. In practice,  $G$  is chosen based on the form of the sparse terminal reward, *e.g.* a max or a sum function. In the learning process, the  $S_\theta$  function is learned by back-propagating errors between the terminal sparse reward and the predicted  $J$  through the aggregation function. In this paper, we use  $\ell_2$  loss that is:

$$\min_{\theta} \frac{1}{2} (J_\theta(\tau) - r(\tau))^2$$

**Spatial Reward Decomposition** An environment usually contains multiple objects. Those objects might re-appear at various spatial locations over time. To further assist the learned knowledge to be transferrable to the new environment, we take advantage of the compositionality of the environment by also decomposing the reward function spatially. More specifically, we divide the contingency-aware observation into smaller sub-regions. For example, in a Cartesian coordinate system, we can divide each coordinate independently and uniformly. With the sub-regions, we re-parametrize the scoring function as  $\sum_{l \in \mathcal{L}} S_\theta(W_l(s_t), \mathbf{a}_t)$ , where  $l$  is the index of the sub-regions and we re-target  $S_\theta$  for the sub-region instead of the whole contingency-aware observation. The intuition is that the smaller sub-regions contains objects or other unit elements that are also building blocks of unseen environments. Such sub-regions become crucial later to generalize to the new environment.

### C. Policy Stage

To solve the novel task with zero interaction, we propose to use planning algorithms to find optimal actions based on the learned scoring function and a learned dynamics model. As shown in the part (c) of Figure 2, we learn a forward dynamics model  $\mathcal{M}_\phi$  based on the exploratory data with a supervised loss function. Specifically, we train a neural network that takes in the action  $\mathbf{a}_t$ , state  $s_t$  and output  $\hat{s}_{t+1}$ , which is an estimate of  $s_{t+1}$ . We use an  $\ell_2$  loss as the objective:

$$\min_{\phi} \frac{1}{2} (\mathcal{M}_\phi(s_t, \mathbf{a}_t) - s_{t+1})^2$$

With the learned dynamics model and the scoring function, we solve an optimization problem using the Model Predictive Control (MPC) algorithm to find the best trajectory for a task  $\mathcal{T}$  in environment  $\mathcal{E}_2$ . The objective of the optimization problem is to minimize  $-J_\theta(\tau')$ . Here we randomly sample multiple action sequences up to length  $H$ , unroll the states based on the current state and the learned dynamics model while computing the cost with the scoring function. We select the action sequence with the minimal cost, and execute the first action in the selected sequence in the environment.

### D. Discussion on zero-shot generalization

Although neural networks have some generalization capabilities, it is still easy to overfit to the training domain. Previous works [39] notice that neural network “memorizes” the optimal action throughout the training process. One can not expect an agent that only remembers optimal actions to generalize well when placed in a new environment. Our method does not suffer from the memorization issue as much as the neural network policy because it can come up with novel solutions, that are not necessarily close to the training examples, since our method learns generalizable model and scores that makes planning possible. The generalization power comes mainly from the smaller building blocks shared across environments as well as the universal dynamics model. This avoids the SAP method to replay the action of the nearest neighbour state in the training data.

Section IV-A.1, Figure 3a 3b compares our method and a neural network RL policy. It shows that in the training environment, RL policy is only slightly worse than our method, however, the RL policy performs much worse than ours in the zero shot generalization case.

## IV. EXPERIMENT

In this section, we study how well the proposed method performs compare to baselines, and the roles of the proposed temporal and spatial reward decompositions. We conduct experiments on two domains: a famous video game “Super Mario Bros” [40] and a robotics blocked reacher task [41].

### A. Experiment on Super Mario Bros

To evaluate our proposed algorithm in a challenging environment, we run our method and baseline methods in the Super Mario Bros environment. This environment features high-dimensional visual observations, which is challenging

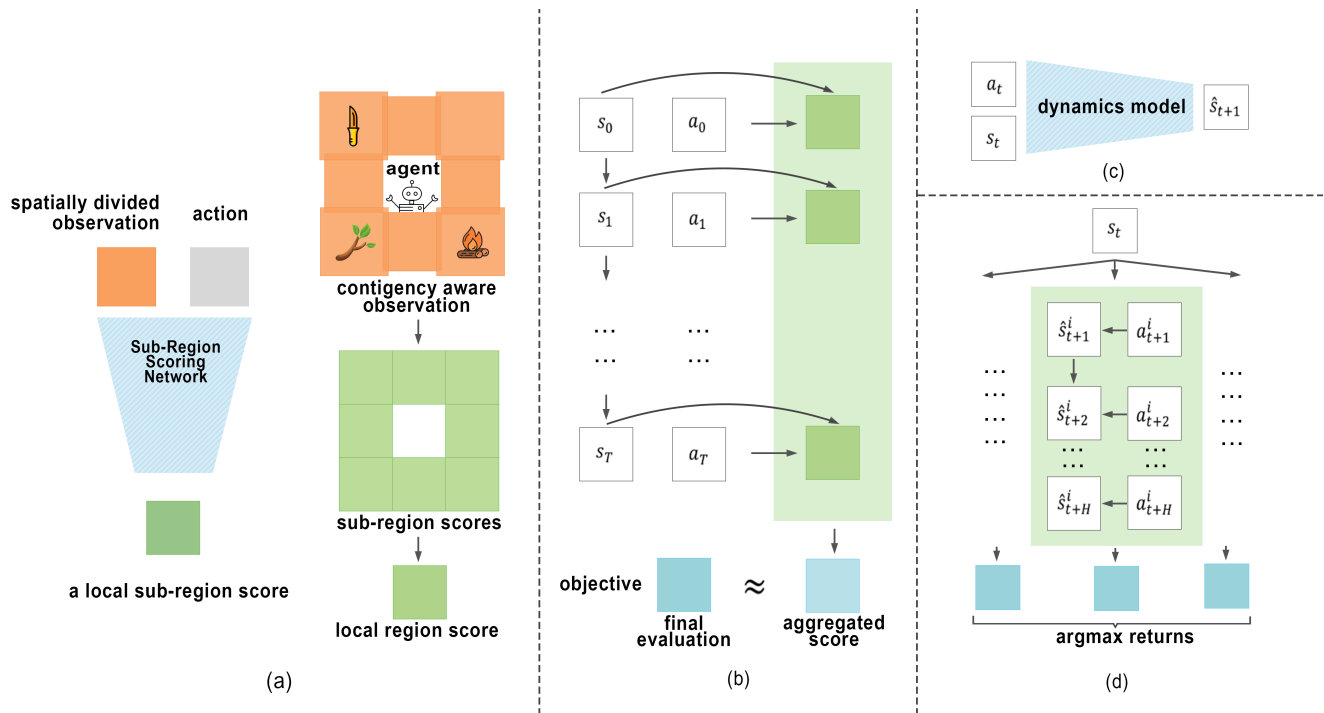


Fig. 2: An overview of the SAP method. (a) For each time step, a scoring network scores contingent sub-regions conditioned on action. (b) we aggregate the prediction over all time steps to fit terminal reward (c)&(d) describe the dynamics learning and planning in mpc in the policy stage.

since we have a large hypothesis space. The original game has  $240 \times 256$  image input and discrete action space with 5 choices. We wrap the environment following [1]. Finally, we obtain a  $84 \times 84$  size 4-frame gray-scale stacked observation. The goal for an agent is to survive and go toward the right as far as possible. We don't have access to the dense reward for each step but the environment returns how far the agent moves towards target at the end of a trajectory as the delayed terminal sparse reward.

1) *Baselines*: We compare our method with various types of baselines, including state-of-art zero-shot RL algorithms, generic policy learning algorithms as well as oracles with much more environment interactions. More details can be found in the appendix.

Exploration Data Exploration Data is the data from which we learn the scores, dynamics model and imitate. The data is collected from noisy sub-optimal version of policy trained using [33]. The average reward on this dataset is a baseline for all other methods.

Behavioral Cloning [42], [43] Behavioral Cloning (BC) learns a mapping from a state to an action on the exploration data using supervised learning. We use cross-entropy loss for predicting the actions.

Model Based with Human Prior Model Based with Human Priors method (MBHP) incorporates model predictive control with predefined human priors, which is +1 score if the agent tries to move or jump toward the right and 0 otherwise.

DARLA [15] DARLA relies on learning a latent state representation that can be transferred from the training envi-

ronments to the testing environment. It achieves this goal by obtaining a disentangled representation of the environment's generative factors before learning to act. We use the latent representation as the observations for a behavioral cloning agent.

RUDDER [32] RUDDER proposes to use LSTM to decompose the delayed sparse reward to dense rewards. It first trains a function  $f(\tau_{1:T})$  predict the terminal reward of a trajectory  $\tau_{1:T}$ . It then use  $f(\tau_{1:t}) - f(\tau_{1:t-1})$  as dense reward at step  $t$  to train RL policies. We change policy training to MPC so this reward can be used in zero-shot setting.

Behavior Clone with Privilege Data Instead of using exploratory trajectories from the environment, we collect a set of near optimal trajectories in the training environment and train a behavior clone agent from it. Note that this is not a fair comparison with other methods, since this method uses better performing training data.

RL curiosity We use a PPO [44] agent that is trained with curiosity driven reward [45] and the final sparse reward in the training environment. This also violates the setting we have as it interacts with the environment. We conduct this experiment to test the generalization ability of an RL agent.

2) *Analysis*: Fig. 3 a and Fig. 3 b show how the above methods performs in the Super Mario Bros. The Explore baseline shows the average training trajectory performance. We found that RUDDER fails to match the demonstration performance. This can attributed to the LSTM in RUDDER not having sufficient supervision signals from the long horizon (2k steps) delayed rewards in the Mario environment. DARLA slightly outperforms the demonstration data in training. Its

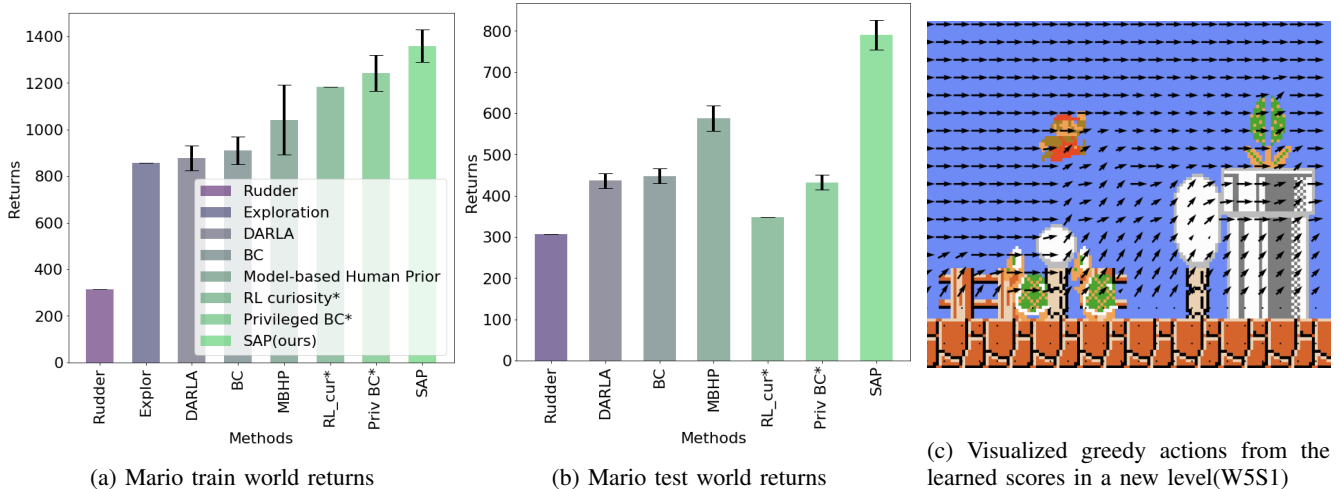


Fig. 3: (a) & (b) The total returns of different methods on the Super Mario Bros. Methods with \* have privilege access to optimal data instead of sub-optimal data. Error bars are shown as 95% confidence interval. See Section IV-A.1 for details.

TABLE I: Ablation of the temporal and spatial reward decomposition.

	W1S1(train)	W2S1(test)
SAP	<b>1359.3</b>	<b>790.1</b>
SAP w/o spatial	1258.0	737.0
SAP w/o spatial temporal	1041.3	587.9

performance is limited by the unsupervised visual disentangled learning step, which is a hard problem in the complex image domain. Behavior Cloning is slightly better than the exploration data, but much worse than our method. MBHP is a model-based approach where human defines the per step reward function. However, it is prohibitive to obtain detailed manual rewards every step. In this task, it is inferior to SAP’s learned priors. We further compare to two methods that unfairly utilize additional privileged information. The Privileged BC method trains on near optimal data, and performs quite well during training; however, it performs much worse during the zero shot testing. The similar trend happens with RL Curiosity, which has the privileged access to online interactions in the training environment.

To conclude, we found generic model free algorithms, including both RL (RL Curiosity), and behavior cloning (Privilege BC) perform well on training data but suffer from severe generalization issue. The zero-shot algorithms (DARLA), BC from exploratory data (BC) and model-based method (MBHP) suffer less from degraded generalization, but they all under-perform our proposed algorithm.

3) *Ablative Studies*: In order to have a more fine-grained understanding of the effect of our newly proposed temporal and spatial reward decomposition, we further conduct ablation studies on those two components. We run two other versions of our algorithm, removing the spatial reward decomposition and the temporal reward decomposition one at a time. More specifically, the SAP w/o spatial method does not divide the observation into sub-regions, but simply use a convolution network to approximate the

scoring function. SAP w/o spatial temporal further removes the temporal reward learning part, and replace the scoring function with a human-designed prior. I.e. it is the same as the MBHP baseline. Please see appendix for more details.

Table I shows the result. We found that both the temporal and the spatial reward learning component contributes to the final superior performance. Without temporal reward decomposition, the agent either have to deal with sparse reward or use a manually specified rewards which might be tedious or impossible to collect. Without the spatial decomposition, the agent might not find correctly which specific region or object is important to the task and hence fail to generalize.

4) *Visualization of learned scores*: In this section, we qualitatively study the induced action by greedily maximizing one-step score. I.e. for any location on the image, we assumes the Super Mario agent were on that location and find the action  $a$  that maximize the learned scoring function  $S_{\theta}(W(s), a)$ . We visualize the computed actions on World 5 Stage 1 (Fig. 3 c) which is visually different from previous tasks. More visualization can be found in <sup>2</sup> In this testing case, we see that the actions are reasonable, such as avoiding obstacles and monsters by jumping over them, even in the face of previously unseen configurations and different backgrounds. However, the “Piranha Plants” are not recognized because all the prior scores are learned from W1S1 where it never appears. More visualization of action maps are available in our videos. Those qualitative studies further demonstrate that the SAP method can assign meaningful scores for different objects in an unsupervised manner. It also produces good actions even in a new environment.

#### B. SAP on the 3-D robotics task

In this section, we further study the SAP method to understand its property with a higher dimensional observation space. We conduct experiments in a 3-D robotics

<sup>2</sup><https://sites.google.com/view/sapnew/home>

environment, BlockedReacher-v0. In this environment, a robot hand is initialized at the left side of a table and tries to reach the right. Between the robot hand and the goal, there are a few blocks standing as obstacles. The task is moving the robot hand to reach a point on  $y = 1.0$  as fast as possible. To test the generalization capability, we create four different configurations of the obstacles, as shown in Figure 4. Figure 4 A is the environment where we collect exploration data from and Figure 4 B, C, D are the testing environments. Note that the exploration data has varying quality, where many of the trajectories are blocked by the obstacles. We introduce more details about this experiment.

1) *Environment.*: In the Blocked Reach environment, we use a 7-DoF robotics arm to reach a specific point. For more details, we refer the readers to [46]. We discretize the robot world into a  $200 \times 200 \times 200$  voxel cube. For the action space, we discretize the actions into two choices for each dimension which are moving 0.5 or -0.5. Hence, in total there are 8 actions. We design four configurations for evaluating different methods as shown in Figure 4. For each configurations, there are three objects are placed in the middle as obstacles.

2) *Applying SAP on the 3D robotic task:* We apply the SAP framework as follows. The original observation is a 25-dimensional continuous state and the action space is a 3-dimensional continuous control. They are discretized into voxels and 8 discrete actions as described in Appendix C.1. The scoring function is a fully-connected neural network that takes in a flattened voxel sub-region. We also train a 3D convolutional neural net as the dynamics model. The dynamics model takes in the contingency-aware observation as well as an action, and outputs the next robot hand location. With the learned scores and the dynamics model, we plan using the MPC method with a horizon of 8 steps.

3) *Architectures for score function and dynamics model.*: For the score function, we train a 1 hidden layer fully-connected neural networks with 128 units. We use ReLu functions as activation except for the last layer. Note that the input 5 by 5 by 5 voxels are flattened before put into the scoring neural network.

For the dynamics model, we train a 3-D convolution neural network that takes in the contingency-aware observation (voxels), action and last three position changes. The voxels contingent to end effector are encoded using three 3d convolution with kernel size 3 and stride 2. Channels of these 3d conv layers are 16, 32, 64, respectively. A 64-unit FC layer is connected to the flattened features after convolution. The action is encoded with one-hot vector connected to a 64-unit FC layer. The last three  $\delta$  positions are also encoded with a 64-unit FC layer. The three encoded features are concatenated and go through a 128-unit hidden FC layer and output predicted change in position. All intermediate layers use ReLu as activation.

4) *Results:* We evaluate similar baselines as in the previous section that is detailed in IV-A.1. In Table II, we compare our method with the MBHP and RUDDER on the 3D robot reaching task. We found that our method needs

TABLE II: Evaluation of SAP, MBHP and Rudder on the 3D Reacher environment. Numbers are the avg. steps to reach the goal. The lower the better. Numbers in the brackets are the 95% confidence interval. ‘‘L’’ denotes the learned dynamics, and ‘‘P’’ denotes the perfect dynamics.

	Config A	Config B	Config C	Config D
SAP(L)	<b>97.53</b> [2.5]	<b>86.53</b> [1.9]	<b>113.3</b> [3.0]	<b>109.2</b> [3.0]
MBHP(L)	124.2[4.5]	102.0[3.0]	160.7[7.9]	155.4[7.2]
Rudder(L)	1852[198]	1901[132]	1933[148]	2000[0.0]
SAP(P)	<b>97.60</b> [2.6]	<b>85.38</b> [1.8]	<b>112.2</b> [3.1]	<b>114.1</b> [3.2]
MBHP(P)	125.3[4.3]	102.4[3.0]	153.4[5.3]	144.9[4.8]
Rudder(P)	213.6[4.2]	194.2[7.0]	208.2[5.6]	201.5[6.3]

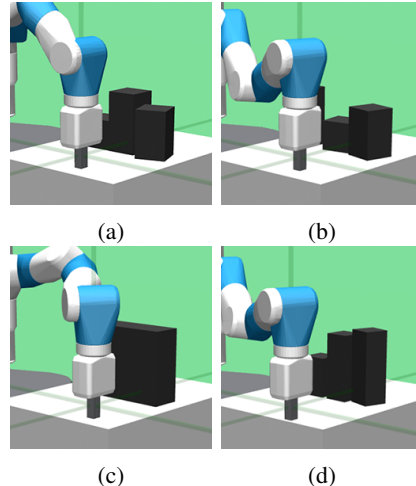


Fig. 4: Four variants of the 3D robot reacher environments. See Section IV-B for details.

significantly fewer steps than the two baselines, in both training environment and testing ones. We find that SAP significantly moves faster to the right because it learns a negative score for crashing into the obstacles. However, the MBHP method, which has +1 positive for each 1 meter moved to the right, would be stuck by the obstacles for a longer duration. When training with RUDDER, the arm also frequently waste time by getting stuck at obstacle. We found that our SAP model is relatively insensitive to the errors in the learned dynamics, such that the performance using the learned dynamics is close to that of perfect dynamics. These experiments show that our method can be applied to robotics environment that can be hard for some algorithms due to the 3-D nature.

## V. CONCLUSION

In this paper, we introduced a new method called SAP that aims to generalize in the new environment without any further interactions by learning the temporally and spatially decomposed rewards. We empirically demonstrate that the newly proposed algorithm outperform previous zero-shot RL method by a large margin, on two challenging environments, i.e. the Super Mario Bros and the 3D Robot Reach. The proposed algorithm along with a wide range of baselines provide a comprehensive understanding of the important aspect zero-shot RL problems.

## VI. ACKNOWLEDGEMENT

The work has been supported by BAIR, BDD, and DARPA XAI. This work has been supported in part by the Zhongguancun Haihua Institute for Frontier Information Technology. Part of the work done while Yang Gao is at UC Berkeley. We also thank Olivia Watkins for discussion.

## REFERENCES

- [1] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu *et al.*, “Human-level control through deep reinforcement learning,” *Nature*, vol. 518, no. 7540, p. 529, 2015.
- [2] J. Schulman, P. Moritz, S. Levine, M. Jordan, and P. Abbeel, “High-dimensional continuous control using generalized advantage estimation,” *arXiv preprint arXiv:1506.02438*, 2015.
- [3] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, “Continuous control with deep reinforcement learning,” *arXiv preprint arXiv:1509.02971*, 2015.
- [4] J. Gibson, *The ecological approach to visual perception*. Psychology Press, 2014.
- [5] R. Dubey, P. Agrawal, D. Pathak, T. L. Griffiths, and A. A. Efros, “Investigating human priors for playing video games,” *arXiv preprint arXiv:1802.10217*, 2018.
- [6] J. Choi, Y. Guo, M. Moczulski, J. Oh, N. Wu, M. Norouzi, and H. Lee, “Contingency-aware exploration in reinforcement learning,” *arXiv preprint arXiv:1811.01483*, 2018.
- [7] D. Q. Mayne, J. B. Rawlings, C. V. Rao, and P. O. Scokaert, “Constrained model predictive control: Stability and optimality,” *Automatica*, vol. 36, no. 6in, pp. 789–814, 2000.
- [8] D. Wang, C. Devin, Q.-Z. Cai, F. Yu, and T. Darrell, “Deep object-centric policies for autonomous driving,” in *ICRA*. IEEE, 2019, pp. 8853–8859.
- [9] E. Jang, C. Devin, V. Vanhoucke, and S. Levine, “Grasp2vec: Learning object representations from self-supervised grasping,” *arXiv preprint arXiv:1811.06964*, 2018.
- [10] C. Devin, P. Abbeel, T. Darrell, and S. Levine, “Deep object-centric representations for generalizable robot learning,” in *ICRA*. IEEE, 2018, pp. 7111–7118.
- [11] G. Zhu, Z. Huang, and C. Zhang, “Object-oriented dynamics predictor,” in *Advances in Neural Information Processing Systems*, 2018, pp. 9804–9815.
- [12] Y. Du and K. Narasimhan, “Task-agnostic dynamics priors for deep reinforcement learning,” *arXiv preprint arXiv:1905.04819*, 2019.
- [13] R. Keramati, J. Whang, P. Cho, and E. Brunskill, “Strategic object oriented reinforcement learning,” *arXiv preprint arXiv:1806.00175*, 2018.
- [14] M. Li, P. Nashikkar, and J. Wang, “Optimizing object-based perception and control by free-energy principle,” *CoRR*, vol. abs/1903.01385, 2019. [Online]. Available: <http://arxiv.org/abs/1903.01385>
- [15] I. Higgins, A. Pal *et al.*, “Darla: Improving zero-shot transfer in reinforcement learning,” in *ICML*. JMLR. org, 2017, pp. 1480–1490.
- [16] S. Sohn, J. Oh, and H. Lee, “Hierarchical reinforcement learning for zero-shot generalization with subtask dependencies,” in *Neurips*, 2018, pp. 7156–7166.
- [17] J. Oh, S. Singh, H. Lee, and P. Kohli, “Zero-shot task generalization with multi-task deep reinforcement learning,” in *ICML*. JMLR. org, 2017, pp. 2661–2670.
- [18] A. Y. Ng, S. J. Russell *et al.*, “Algorithms for inverse reinforcement learning,” in *ICML*, 2000.
- [19] K. Bogert and P. Doshi, “Multi-robot inverse reinforcement learning under occlusion with state transition estimation,” in *AAMAS*, 2015, pp. 1837–1838.
- [20] S. Wang, R. Rosenfeld, Y. Zhao, and D. Schuurmans, “The latent maximum entropy principle,” in *Proceedings IEEE International Symposium on Information Theory*. IEEE, 2002, p. 131.
- [21] K. Bogert, J. F.-S. Lin, P. Doshi, and D. Kulic, “Expectation-maximization for inverse reinforcement learning with hidden data,” in *AAMAS*, 2016, pp. 1034–1042.
- [22] J. Choi and K.-E. Kim, “Inverse reinforcement learning in partially observable environments,” *Journal of Machine Learning Research*, vol. 12, no. Mar, pp. 691–730, 2011.
- [23] U. Syed and R. E. Schapire, “A game-theoretic approach to apprenticeship learning,” in *Advances in neural information processing systems*, 2008, pp. 1449–1456.
- [24] S. Levine and P. Abbeel, “Learning neural network policies with guided policy search under unknown dynamics,” in *NIPS*, 2014, pp. 1071–1079.
- [25] J. Bagnell, J. Chestnutt, D. M. Bradley, and N. D. Ratliff, “Boosting structured prediction for imitation learning,” in *Advances in Neural Information Processing Systems*, 2007, p. 1153.
- [26] A. Y. Ng, “Feature selection,  $l_1$  vs.  $l_2$  regularization, and rotational invariance,” in *Proceedings of the twenty-first international conference on Machine learning*. ACM, 2004, p. 78.
- [27] X. Xie, C. Li, C. Zhang, Y. Zhu, and S.-C. Zhu, “Learning virtual grasp with failed demonstrations via bayesian inverse reinforcement learning,” in *IROS*, 2019.
- [28] A. Y. Ng, D. Harada, and S. Russell, “Policy invariance under reward transformations: Theory and application to reward shaping,” in *ICML*, vol. 99, 1999, pp. 278–287.
- [29] B. Marthi, “Automatic shaping and decomposition of reward functions,” in *Proceedings of the 24th International Conference on Machine learning*. ACM, 2007, pp. 601–608.
- [30] M. Grzes̄ and D. Kudenko, “Online learning of shaping rewards in reinforcement learning,” *Neural Networks*, vol. 23, no. 4, pp. 541–550, 2010.
- [31] M. Marashi, A. Khalilian, and M. E. Shiri, “Automatic reward shaping in reinforcement learning using graph analysis,” in *ICCKE*. IEEE, 2012, pp. 111–116.
- [32] J. A. Arjona-Medina, M. Gillhofer, M. Widrich, T. Unterthiner, J. Brandstetter, and S. Hochreiter, “Rudder: Return decomposition for delayed rewards,” *arXiv preprint arXiv:1806.07857*, 2018.
- [33] D. Pathak, P. Agrawal, A. A. Efros, and T. Darrell, “Curiosity-driven exploration by self-supervised prediction,” in *CVPR workshops*, 2017, pp. 16–17.
- [34] J. Schmidhuber, “A possibility for implementing curiosity and boredom in model-building neural controllers,” in *Proc. of the international conference on simulation of adaptive behavior: From animals to animats*, 1991, p. 222–227.
- [35] H. Tang, R. Houthoofd *et al.*, “# exploration: A study of count-based exploration for deep reinforcement learning,” in *Neurips*, 2017, pp. 2753–2762.
- [36] A. L. Strehl and M. L. Littman, “An analysis of model-based interval estimation for markov decision processes,” *Journal of Computer and System Sciences*, vol. 74, no. 8, pp. 1309–1331, 2008.
- [37] M. Andrychowicz, F. Wolski, A. Ray, J. Schneider, R. Fong, P. Welinder, B. McGrew, J. Tobin, P. Abbeel, and W. Zaremba, “Hindsight experience replay,” in *Neurips*, 2017, pp. 5048–5058.
- [38] A. Dosovitskiy and V. Koltun, “Learning to act by predicting the future,” *arXiv preprint arXiv:1611.01779*, 2016.
- [39] C. Packer, K. Gao, J. Kos, P. Kr̄ahenb̄uhl, V. Koltun, and D. Song, “Assessing generalization in deep reinforcement learning,” *arXiv preprint arXiv:1810.12282*, 2018.
- [40] C. Kauten, “Super Mario Bros for OpenAI Gym,” GitHub, 2018. [Online]. Available: <https://github.com/Kautenja/gym-super-mario-bros>
- [41] Z. Huang, F. Liu, and H. Su, “Mapping state space using landmarks for universal goal reaching,” *arXiv preprint arXiv:1908.05451*, 2019.
- [42] D. A. Pomerleau, “Alvinn: An autonomous land vehicle in a neural network,” in *Advances in neural information processing systems*, 1989, pp. 305–313.
- [43] M. Bain and C. Sammut, “A framework for behavioural cloning,” in *Machine Intelligence 15*, 1995, pp. 103–129.
- [44] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, “Proximal policy optimization algorithms,” *arXiv preprint arXiv:1707.06347*, 2017.
- [45] Y. Burda, H. Edwards, A. Storkey, and O. Klimov, “Exploration by random network distillation,” *arXiv preprint arXiv:1810.12894*, 2018.
- [46] M. Plappert, M. Andrychowicz *et al.*, “Multi-goal reinforcement learning: Challenging robotics environments and request for research,” *arXiv preprint arXiv:1802.09464*, 2018.
- [47] W. Lotter, G. Kreiman, and D. Cox, “Deep predictive coding networks for video prediction and unsupervised learning,” *arXiv preprint arXiv:1605.08104*, 2016.
- [48] C. Finn, I. Goodfellow, and S. Levine, “Unsupervised learning for physical interaction through video prediction,” in *Advances in neural information processing systems*, 2016, pp. 64–72.

### A. Hidden Reward Gridworld

1) *Environment.*: In the gridworld environment, each entry correspond to a feature vector with noise based on the type of object in it. Each feature is a length 16 vector whose entries are uniformly sampled from  $[0, 1]$ . Upon each feature, we add a small random noise from a normal distribution with  $\mu = 0, \sigma = 0.05$ . The outer-most entries correspond to padding objects whose rewards are 0. The action space includes move toward four directions up, down, left, right. If an agent attempts to take an action which leads to outside of our grid, it will be ignored by the environment.

2) *Architectures for score function and dynamics model.*: We train a two layer fully connected neural networks with 32 and 16 hidden units respectively and a ReLU activation function to approximate the score for each grid.

In this environment, we do not have a learned dynamics model.

**Hyperparameters.** During training, we use an adam optimizer with learning rate  $1e-3$ ,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ . The learning rate is reduced to  $1e-4$  after 30000 iterations. The batchsize is 128. We use  $horizon = 4$  as our planning horizon.

### B. Super Mario Bros

1) *Environment.*: We wrap the original Super Mario environments with additional wrappers. We wrap the action space into 5 discrete joystick actions, none, walk right, jump right, run right and hyper jump right. We follow [45] to add a sticky action wrapper that repeats the last action with a probability of 20%. Besides this, we follow add the standard wrapper as in past work [1].

2) *Applying SAP on Super Mario Bros:* We apply our SAP framework as follows. We first divide the egocentric region around the agent into eight 12 by 12 pixel sub-regions based on relative position as illustrated in Figure. 5 in the Appendix. Each sub-region is scored by a CNN, which has a final FC layer to output a score matrix. The matrix has the shape  $\text{dim}(\text{action}) \times \text{dim}(\text{relative position})$ , which are 5 and 8 respectively. Then an action selector and a sub-region selector jointly select row corresponding to the agent’s action and the column corresponding to the relative position. The sum of all the sub-region scores forms the egocentric region score. Then we minimize the  $\ell_2$  loss between the aggregated egocentric region scores along the trajectory and the terminal reward. In addition, we add an regularization term to the loss to achieve better stability for long horizon. We take the  $\ell_1$  norm of score matrices across trajectory to add to loss term. This encourages the scoring function to center predicted loss terms around 0 and enforce numerical sparsity. A dynamics model is also learned by training another CNN. The dynamics model takes in a 30 by 30 size crop around the agent, the agent’s location as well a one-hot action vector. Instead of outputting a full generated image, we only predict the future location of the agent recursively. We avoid video predictive models because it suffers the blurry effect when

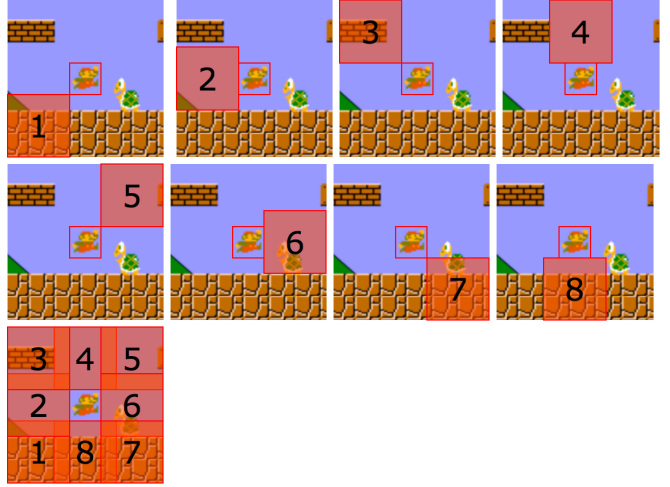


Fig. 5: A visualization of the sub-regions in the Super Mario Bros game. In this game, there are in total 8 sub-regions.

predicting long term future [47], [48]. We plan with the learned scores and dynamics model with a standard MPC algorithm with random actions that looks ahead 10 steps.

3) *Architectures for score function and dynamics model.*: For the score function, we train a CNN taking each 12px by 12px sub-region as input with 2 conv layers and 1 hidden fully connected layers. For each conv layer, we use a filter of size 3 by 3 with stride 2 with number of output channels equals to 8 and 16 respectively. “Same padding” is used for each conv layer. The fully connected layers have 128 units. Relu functions are applied as activation except the last layer.

For the dynamics model, we train a neural network with the following inputs: a. 30 by 30 egocentric observation around mario. b. current action along with 3 recent actions encoded in one-hot tensor. c. 3 most recent position shifts. d. one-hot encoding of the current planning step. Input a is encoded with 4 sequential conv layers with kernel size 3 and stride 2. Output channels are 8, 16, 32, 64 respectively. A global max pooling follows the conv layers. Input b, c, d are each encoded with a 64 node fc layer. The encoded results are then concatenated and go through a 128 units hidden fc layer. This layer connects to two output heads, one predicting shift in location and one predicting “done” with sigmoid activation. Relu function is applied as activation for all intermediate layers.

4) *Hyperparameters.*: During training, we use an adam optimizer with learning rate  $3e-4$ ,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ . The batchsize is 256 for score function training and 64 for dynamics model. We use  $horizon = 10$  as our planning horizon. We use a discount factor  $\gamma = 0.95$  and 128 environments in our MPC. We use  $1e-4$  for the regularization term for matrices entries.

5) *More on training:* In the scoring function training, each data point is a tuple of a down sampled trajectory and a calculated score. We down sample the trajectory in the exploration data by taking data from every two steps. Half the the trajectories ends with a “done”(death) event and half



are not. For those ends with “done”, the score is the distance mario traveled by mario at the end. For the other trajectories, the score is the distance mario traveled by the end plus a mean future score. The mean future score of a trajectory is defined to be the average extra distance traveled by longer (in terms of distance) trajectories than our trajectory. We note that all the information are contained in the exploration data.

6) *More Details on Baselines: Exploration Data* Exploration Data is the data from which we learn the scores, dynamics model and imitate. The data is collected from a suboptimal policy described in Appendix B.6. The average reward on this dataset is a baseline for all other methods. This is omitted in new tasks because we only know the performance in the environment where the data is collected. We train a policy with only curiosity as rewards [33]. However, we early stopped the training after  $5e7$  steps which is far from the convergence at  $1e9$  steps. We further added an  $\epsilon$ -greedy noise when sampling demonstrations with  $\epsilon = 0.4$  for 20000 episodes and  $\epsilon = 0.2$  for 10000 episodes.

**Behavioral Cloning [42], [43]** Behavioral Cloning (BC) learns a mapping from a state to an action on the exploration data using supervised learning. We use cross-entropy loss for predicting the actions.

**Model Based with Human Prior** Model Based with Human Priors method (MBHP) incorporates model predictive control with predefined human priors, which is +1 score if the agent tries to move or jump toward the right and 0 otherwise. MBHP replaces the scoring-aggregation step of our method by a manually defined prior. We note that it might be hard to design human priors in other tasks. As super mario is a deterministic environment, we noticed pure behavior cloning trivially get stuck at a tube at the very beginning of level 1-1 and die at an early stage of 2-1. Thus we select action using sampling from output logits instead of taking argmax.

**DARLA [15]** DARLA relies on learning a latent state representation that can be transferred from the training environments to the testing environment. It achieves this goal by obtaining a disentangled representation of the environment’s generative factors before learning to act. We use the latent representation as the observations for a behavioral cloning agent. We re-implemented and fine-tuned a disentangled representation of observation space as described in [15] for mario environment. We used 128 latent dimensions in both DAE and beta-VAE with  $\beta = 0.1$ . The VAE is trained with batch size 64, learning rate  $1e-4$ . The DAE is trained with batch size 64 and learning rate  $1e-3$ . We tune until the test visualization on training environment is good. We then train a behavior cloning on the learned disentangled representation to benchmark on both training and testing environment.

**RUDDER [32]** RUDDER proposes to use LSTM to decompose the delayed sparse reward to dense rewards. It first trains a function  $f(\tau_{1:T})$  predict the terminal reward of a trajectory  $\tau_{1:T}$ . It then use  $f(\tau_{1:t}) - f(\tau_{1:t-1})$  as dense reward at step  $t$  to train RL policies. We change policy training to MPC so this reward can be used in zero-shot setting. We re-implemented [32] on both mario and robot enviroment and fine-tuned it respectively. We down-sample mario trajectory

TABLE III: Ablation Study for number of planning steps in MPC based methods. The averaged return is reported.

	World 1 Stage 1		
	plan 8	plan 10	plan 12
SAP	<b>1341.5</b>	<b>1359.3</b>	<b>1333.5</b>
MBHP	1193.8	1041.3	1112.1
	World 2 Stage 1(new task)		
	plan 8	plan 10	plan 12
SAP	<b>724.2</b>	<b>790.1</b>	<b>682.6</b>
MBHP	546.4	587.9	463.8



Fig. 6: More visualizations on the greedy action map on W1S1, W2S1(new task) and W5S1(new task). Note the actions can be different from the policy from MPC.

by a factor of 2 to feed into a special LSTM in RUDDER’s source code. For mario environment, the feature vector for each time step is derived from a convolution network with channels 1, 3, 8, 16, 32; kernal sizes 3, 3, 3, 3 respectively, followed by a fc layer with 64 output features. The 64 is then feed into each LSTM step. For robot reaching environment, we directly use a fc layer to down sample observation feature from 3250 to 64 and feed into LSTM. Both models are trained with batch size 256 under learning rate  $2e-4$ , with weight decay of  $5e-3$ .

**Behavior Clone with Privilege Data** Instead of using exploratory trajectories from the environment, we collect a set of near optimal trajectories in the training environment and train a behavior clone agent from it. Note that this is not a fair comparison with other methods, since this method uses better performing training data.

**RL curiosity** We use a PPO [44] agent that is trained with curiosity driven reward [45] and the final sparse reward in the training environment. We limit the training steps to 10M. This is also violating the setting we have as it interacts with the environment. We conduct this experiment to test the generalization ability of a RL agent.

7) *Additional Ablations: Ablation of Planning Steps* In this section, we conduct additional ablative experiments to evaluate the effect of the planning horizon in a MPC method. In Table. III, we see that our method fluctuates a little with different planning steps in a relatively small range and outperforms baselines constantly. In the main paper, we choose  $horizon = 10$ . We find that when plan steps are larger such as 12, the performance does not improve monotonically. This might be due to the difficult to predict long range future with a learned dynamics model.

8) *Additional visualization:* In this section, we present additional visualization for qualitative study. In Figure. 6, we

see that on a few randomly sampled frames, even the greedy action can be meaningful for most of the cases. We see the agent intend to jump over obstacles and avoid dangerous monsters.

In Figure. 7, we show the scores of a given state-action pair and find that the scores fulfill the human prior. For example, in Figure. 7a, we synthetically put the mario agent in 8 relative positions to “koopas” conditioned on the action “move right”. The score is significantly lower when the agent’s position is to the left of “koopas” compared to other position. In Figure. 7b, it is the same setup as in Figure. 7a but conditioned on the action “jump”. We find that across Figure. 7a and Figure. 7b the left position score of Figure. 7b is smaller than that of Figure. 7a which is consistent with human priors. In Figure. 7c and Figure. 7c, we substitute the object from “koopas” to the ground. We find that on both Figure. 7c and Figure. 7c the score are similar for the top position which means there is not much difference between different actions.

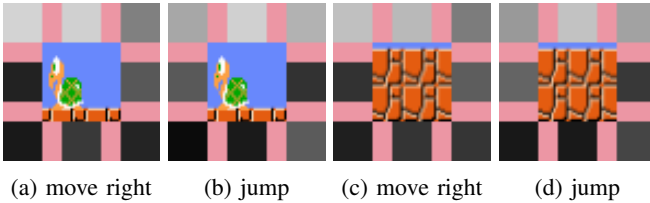
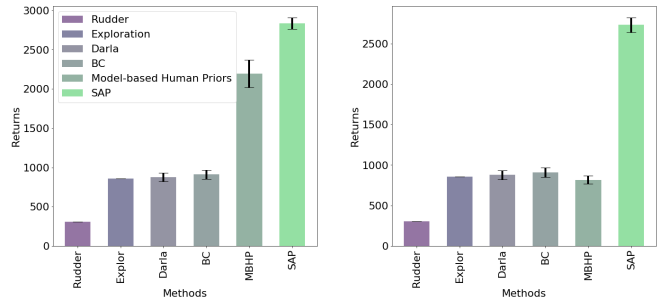


Fig. 7: Visualization of the learned score on pre-extracted objects. The grayscale area is the visualized score and the pink area is a separator. Best viewed in color. In (a), we synthetically put the mario agent in 8 relative positions to “koopas” conditioned on the action “move right”. The score is significantly lower when the agent’s position is to the left of “koopas” compared to other position. In (b), it is the same setup as in (a) but conditioned on the action “jump”. We find that across (a) and (b) the left position score of (b) is smaller than that of (a) which is consistent with human priors. In (c) and (d), we substitute the object from “koopas” to the ground. We find that on both (c) and (d) the score are similar for the top position which means there is not much difference between different actions. Note this figure is only for visualizations and we even put the agent in positions that cannot be achieved in the actual game.

### 9) Additional Results with Ground Truth Dynamics Model and No Done Signal:

#### C. Robotics Blocked Reach

1) *Environment.*: In the Blocked Reach environment, a 7-DoF robotics arm is manipulated for a specific task. For more details, we refer the readers to [46]. We discretize the robot world into a  $200 \times 200 \times 200$  voxel cube. For the action space, we discretize the actions into two choices for each dimension which are moving 0.5 or -0.5. Hence, in total there are 8 actions. We design four configurations for evaluating different methods as shown in Figure 4. For each configurations, there are three objects are placed in



(a) GT dynamics model on WIS1 (train) (b) GT dynamics model & no done on WIS1 (train)

Fig. 8: Training Performance with perfect dynamics and no done signal.

the middle as obstacles. The height of the objects in each configuration are (0.05, 0.1, 0.08), (0.1, 0.05, 0.08), (0.12, 1.12, 0.12), (0.07, 0.11, 0.12).

2) *Applying SAP on the 3D robotic task:* We apply the SAP framework as follows. The original observation is a 25-dimensional continuous state and the action space is a 3-dimensional continuous control. They are discretized into voxels and 8 discrete actions as described in Appendix C.1. In this environment, the egocentric region is set to a  $15 \times 15 \times 15$  cube of voxels around robot hand end effector. We divide this cube into  $27 \times 5 \times 5 \times 5$  sub-regions. The scoring function is a fully-connected neural network that takes in a flattened voxel sub-region and outputs the score matrix with a shape of  $26 \times 8$ . The scores for each step are aggregated by a sum operator along the trajectory. We also train a 3D convolutional neural net as the dynamics model. The dynamics model takes in a  $15 \times 15 \times 15$  egocentric region as well as an action, and outputs the next robot hand location. With the learned scores and the dynamics model, we plan using the MPC method with a horizon of 8 steps.

3) *Architectures for score function and dynamics model.*: For the score function, we train a 1 hidden layer fully-connected neural networks with 128 units. We use Relu functions as activation except for the last layer. Note that the input 5 by 5 by 5 voxels are flattened before put into the scoring neural network.

For the dynamics model, we train a 3-D convolution neural network that takes in a egocentric region (voxels), action and last three position changes. The 15 by 15 by 15 egocentric voxels are encoded using three 3d convolution with kernel size 3 and stride 2. Channels of these 3d conv layers are 16, 32, 64, respectively. A 64-unit FC layer is connected to the flattened features after convolution. The action is encoded with one-hot vector connected to a 64-unit FC layer. The last three  $\delta$  positions are also encoded with a 64-unit FC layer. The three encoded features are concatenated and go through a 128-unit hidden FC layer and output predicted change in position. All intermediate layers use relu as activation.

4) *Hyperparameters.*: During training, we use an adam optimizer with learning rate  $3e-4$ ,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ . The batchsize is 128 for score function training and 64 for

dynamics model. We use  $horizon = 8$  as our planning horizon. We use  $2e-5$  as the weight for matrices entry regularization.

5) *Baselines*: Our model based human prior baseline in the blocked robot environment is a 8-step MPC where score for each step is the y component of the action vector at that step.

We omit the Behavioral cloning baselines, which imitates exploration data, as a consequence of two previous results.