

Online Abstraction with MDP Homomorphisms for Deep Learning*

Ondrej Biza
Czech Technical University
Prague, Czech Republic
bizaondr@fit.cvut.cz

Robert Platt
Northeastern University
Boston, MA, USA
rplatt@ccs.neu.edu

ABSTRACT

Abstraction of Markov Decision Processes is a useful tool for solving complex problems, as it can ignore unimportant aspects of an environment, simplifying the process of learning an optimal policy. In this paper, we propose a new algorithm for finding abstract MDPs in environments with continuous state spaces. It is based on MDP homomorphisms, a structure-preserving mapping between MDPs. We demonstrate our algorithm's ability to learn abstractions from collected experience and show how to reuse the abstractions to guide exploration in new tasks the agent encounters. Our novel task transfer method outperforms baselines based on a deep Q-network in the majority of our experiments. The source code is at https://github.com/ondrejba/aamas_19.

KEYWORDS

reinforcement learning; abstraction; mdp homomorphism; transfer learning; deep learning

ACM Reference Format:

Ondrej Biza and Robert Platt. 2019. Online Abstraction with MDP Homomorphisms for Deep Learning. In *Proc. of the 18th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2019), Montreal, Canada, May 13-17, 2019*, IFAAMAS, 9 pages.

1 INTRODUCTION

The ability to create useful abstractions automatically is a critical tool for an autonomous agent. Without this, the agent is condemned to plan or learn policies at a relatively low level of abstraction, and it becomes hard to solve complex tasks. What we would like is the ability for the agent to learn new skills or abstractions over time that gradually increase its ability to solve challenging tasks. This paper explores this in the context of reinforcement learning.

There are two main approaches to abstraction in reinforcement learning: temporal abstraction and state abstraction. In temporal abstraction, the agent learns multi-step skills, i.e. policies for achieving subtasks. In state abstraction, the agent learns to group similar states together for the purposes of decision making. For example, for handwriting a note, it may be irrelevant whether the agent is

*This research was partially supported by grant no. GA18-18080S of the Grant Agency of the Czech Republic, grant no. EF15_003/0000421 of the Ministry of Education, Youth and Sports of the Czech Republic and by the National Science Foundation through IIS-1427081, IIS-1724191, and IIS-1724257, NASA through NNX16AC48A and NNX13AQ85G, ONR through N000141410047, Amazon through an ARA to Platt, and Google through a FRA to Platt. We thank Tokine Atsuta, Braylan Impata and Nao Ouyang for useful feedback on the manuscript.

Proc. of the 18th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2019), N. Agmon, M. E. Taylor, E. Elkind, M. Veloso (eds.), May 13-17, 2019, Montreal, Canada. © 2019 International Foundation for Autonomous Agents and Multiagent Systems (www.ifaamas.org). All rights reserved.

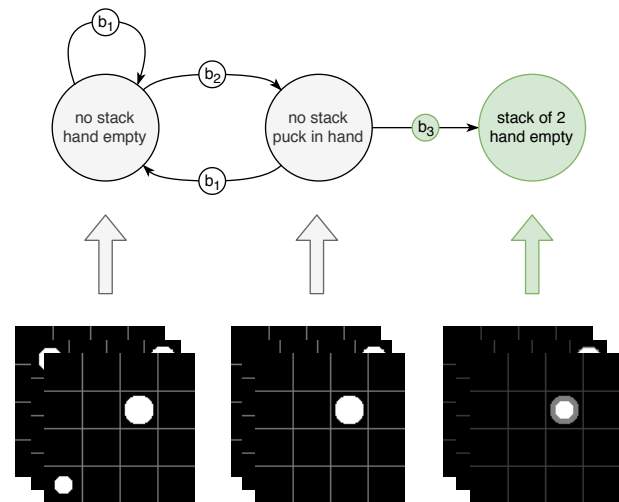


Figure 1: Abstraction for the task of stacking two pucks on top of one another. The diagram shows a minimal quotient MDP (top) that is homomorphic to the underlying MDP (bottom). The minimal MDP has three states, the last of them being the goal state, and four actions. Each action is annotated with the state-action block that induced it. Two actions are annotated with b_1 because they both lead to the first state.

holding a pencil or a pen. In the context of the Markov Decision Process (MDP), state abstraction can be understood using an elegant approach known as the MDP homomorphism framework [13]. An MDP homomorphism is a mapping from the original MDP to a more compact MDP that preserves the important transition and reward structure of the original system. Given an MDP homomorphism to a compact MDP, one may solve the original problem by solving the compact MDP and then projecting those solutions back onto the original problem. Figure 1 illustrates this in the context of a toy-domain puck stacking problem. The bottom left of Figure 1 shows two pucks on a 4×4 grid. The agent must pick up one of the pucks (bottom middle of Figure 1) and place it on top of the other puck (bottom right of Figure 1). The key observation to make here is that although there are many different two-puck configurations (bottom right of Figure 1), they are all equivalent in the sense that the next step is for the agent to pick up one of the pucks. In fact, for puck stacking, the entire system can be summarized by the three-state MDP shown at the top of Figure 1. This compact MDP is clearly a useful abstraction for this problem.

Although MDP homomorphisms are a useful mechanism for abstraction, it is not yet clear how to learn the MDP homomorphism

mapping from experience in a model-free scenario. This is particularly true for a deep reinforcement learning context where the state space is effectively continuous. The closest piece of related work is probably that of [21] who study the MDP homomorphism learning problem in a narrow context. This paper considers the problem of learning general MDP homomorphisms from experience. We make the following key contributions:

#1: We propose an algorithm for learning MDPs homomorphisms from experience in both discrete and continuous state spaces (Subsection 4.2). The algorithm groups together state-action pairs with similar behaviors, creating a partition of the state-action space. The partition then induces an abstract MDP homomorphic to the original MDP. We prove the correctness of our method in Section 5.
#2: Our abstraction algorithm requires a learning component. We develop a classifier based on the Dilated Residual Network [22] that enables our algorithm to handle medium-sized environments with continuous state spaces. We include several augmentations, such as oversampling the minority classes and thresholding the confidences of the predictions (Subsection 4.3). We test our algorithm in two environments: a set of continuous state space puck stacking tasks, which leverage the convolutional network (Subsection 6.1), and a discrete state space blocks world task (Subsection 6.2), which we solve with a decision tree.

#3: We propose a transfer learning method for guiding exploration in a new task with a previously learned abstract MDP (Subsection 4.4). Our method is based on the framework of options [18]: it can augment any existing reinforcement learning agent with a new set of temporally-extended actions. The method outperforms two baselines based on a deep Q-network [9] in the majority of our experiments.

2 BACKGROUND

2.1 Reinforcement Learning

An agent's interaction with an environment can be modeled as a Markov Decision Process (MDP, [1]). An MDP is a tuple $\langle S, A, \Phi, P, R \rangle$, where S is the set of states, A is the set of actions, $\Phi \subset S \times A$ is the state-action space (the set of available actions for each state), $P(s, a, s')$ is the transition function and $R(s, a)$ is the reward function.

We use the framework of options [18] to transfer knowledge between similar tasks. An option $\langle I, \pi, \beta \rangle$ is a temporally extended action: it can be executed from the set of states I and selects primitive actions with a policy π until it terminates. The probability of terminating in each state is expressed by $\beta : S \rightarrow [0, 1]$.

2.2 Abstraction with MDP homomorphisms

Abstraction in our paper aims to group similar state-action pairs from the state-action space Φ . The grouping can be described as a partitioning of Φ .

Definition 2.1. A partition of an MDP $M = \langle S, A, \Phi, P, R \rangle$ is a partition of Φ . Given a partition B of M , the *block transition probability* of M is the function $T : \Phi \times B|S \rightarrow [0, 1]$ defined by $T(s, a, [s']_{B|S}) = \sum_{s'' \in [s']_{B|S}} P(s, a, s'')$.

Definition 2.2. A partition B' is a *refinement* of a partition B , $B' \ll B$, if and only if each block of B' is a subset of some block of B .

The partition of Φ is projected on the state space S to obtain a grouping of states.

Definition 2.3. Let B be a partition of $Z \subseteq X \times Y$, where X and Y are arbitrary sets. For any $x \in X$, let $B(x)$ denote the set of distinct blocks of B containing pairs of which x is a component, that is, $B(x) = \{[(w, y)]_B \mid (w, y) \in Z, w = x\}$. The *projection of B onto X* is the partition $B|X$ of X such that for any $x, x' \in X$, $[x]_{B|X} = [x']_{B|X}$ if and only if $B(x) = B(x')$.

Next, we define two desirable properties of a partition over Φ .

Definition 2.4. A partition B of an MDP $M = \langle S, A, \Phi, P, R \rangle$ is said to be *reward respecting* if $(s_1, a_1) \equiv_B (s_2, a_2)$ implies $R(s_1, a_1) = R(s_2, a_2)$ for all $(s_1, a_1), (s_2, a_2) \in \Phi$.

Definition 2.5. A partition B of an MDP $M = \langle S, A, \Phi, P, R \rangle$ has the *stochastic substitution property* (SSP) if for all $(s_1, a_1), (s_2, a_2) \in \Phi$, $(s_1, a_1) \equiv_B (s_2, a_2)$ implies $T(s_1, a_1, [s]_{B|S}) = T(s_2, a_2, [s]_{B|S})$ for all $[s]_{B|S} \in B|S$.

Having a partition with these properties, we can construct the *quotient MDP* (we also call it the *abstract MDP*).

Definition 2.6. Given a reward respecting SSP partition B of an MDP $M = \langle S, A, \Phi, P, R \rangle$, the *quotient MDP M/B* is the MDP $\langle S', A', \Phi', P', R' \rangle$, where $S' = B|S$; $A' = \bigcup_{[s]_{B|S} \in S'} A'_{[s]_{B|S}}$ where $A'_{[s]_{B|S}} = \{a'_1, a'_2, \dots, a'_{\eta(s)}\}$ for each $[s]_{B|S} \in S'$; P' is given by $P'([s]_f, a'_i, [s']_f) = T_b([(s, a_i)]_B, [s']_{B|S})$ and R' is given by $R'([s]_{B|S}, a'_i) = R(s, a_i)$. $\eta(s)$ is the number of distinct classes of B that contain a state-action pair with s as the state component.

We want the quotient MDP to retain the structure of the original MDP while abstracting away unnecessary information. MDP homomorphism formalizes this intuition.

Definition 2.7. An *MDP homomorphism* from $M = \langle S, A, \Phi, P, R \rangle$ to $M' = \langle S', A', \Phi', P', R' \rangle$ is a tuple of surjections $\langle f, \{g_s : s \in S\} \rangle$ with $h(s, a) = (f(s), g_s(a))$, where $f : S \rightarrow S'$ and $g_s : A \rightarrow A'$ such that $R(s, a) = R'(f(s), g_s(a))$ and $P(s, a, f^{-1}(f(s'))) = P'(f(s), g_s(a), f(s'))$. We call M' a *homomorphic image* of M under h .

The following theorem states that the quotient MDP defined above retains the structure of the original MDP.

THEOREM 2.8 ([13]). *Let B be a reward respecting SSP partition of MDP $M = \langle S, A, \Phi, P, R \rangle$. The quotient MDP M/B is a homomorphic image of M .*

Computing the optimal state-action value function in the quotient MDP usually requires fewer computations, but does it help us act in the underlying MDP? The last theorem states that the optimal state-action value function *lifted* from the minimized MDP is still optimal in the original MDP:

THEOREM 2.9 (OPTIMAL VALUE EQUIVALENCE, [13]). *Let $M' = \langle S', A', \Phi', P', R' \rangle$ be the homomorphic image of the MDP $M = \langle S, A, \Phi, P, R \rangle$ under the MDP homomorphism $h(s, a) = (f(s), g_s(a))$. For any $(s, a) \in \Phi$, $Q^*(s, a) = Q^*(f(s), g_s(a))$.*

Algorithm 1 Abstraction

```

1: procedure ABSTRACTION
2:    $E \leftarrow$  collect initial experience with an arbitrary policy  $\pi$ 
3:    $g \leftarrow$  a classifier for state-action pairs
4:    $B \leftarrow \text{OnlinePartitionIteration}(E, g)$ 
5:    $M' \leftarrow$  a quotient MDP constructed from  $B$  according to
      Definition 2.6
6: end procedure

```

3 RELATED WORK

Balaraman Ravindran proposed Markov Decision Process (MDP) homomorphism together with a sketch of an algorithm for finding homomorphisms (i.e. finding the minimal MDP homomorphic to the underlying MDP) given the full specification of the MDP in his Ph.D. thesis [13]. The first and only algorithm (to the best of our knowledge) for finding homomorphisms from experience (online) [21] operates over Controlled Markov Processes (CMP), an MDP extended with an output function that provides more supervision than the reward function alone. Homomorphisms over CMPs were also used in [20] to find objects that react the same to a defined set of actions.

An approximate MDP homomorphism [14] allows aggregating state-action pairs with similar, but not the same dynamics. It is essential when learning homomorphisms from experience in non-deterministic environments because the estimated transition probabilities for individual state-action pairs will rarely be the same, which is required by the MDP homomorphism. Taylor et al. [19] built upon this framework by introducing a similarity metric for state-action pairs as well as an algorithm for finding approximate homomorphisms.

Sorg et al. [16] developed a method based on homomorphisms for transferring a predefined optimal policy to a similar task. However, their approach maps only states and not actions, requiring actions to behave the same across all MDPs. Soni et al. and Rajendran et al. [12, 15] also studied skill transfer in the framework of MDP homomorphisms. Their works focus on the problem of transferring policies between discrete or factored MDPs with pre-defined mappings, whereas our primary contribution is the abstraction of MDPs with continuous state spaces.

4 METHODS

We solve the problem of abstracting an MDP with a discrete or continuous state-space and a discrete action space. The MDP can have an arbitrary reward function, but we restrict the transition function to be deterministic. This restriction simplifies our algorithm and makes it more sample-efficient (because we do not have to estimate the transition probabilities for each state-action pair).

This section starts with an overview of our abstraction process (Subsection 1), followed by a description of our algorithm for finding MDP homomorphisms (Subsection 4.2). We describe several augmentations to the base algorithm that increase its robustness in Subsection 4.3. Finally, Subsection 4.4 contains the description of our transfer learning method that leverages the learned MDP homomorphism to speed up the learning of new tasks.

Algorithm 2 Online Partition Iteration

```

Input: Experience  $E$ , classifier  $g$ .
Output: Reward respecting SSP partition  $B$ .
1: procedure ONLINEPARTITIONITERATION
2:    $B \leftarrow \{E\}, B' \leftarrow \{\}$ 
3:    $B \leftarrow \text{SplitRewards}(B)$ 
4:   while  $B \neq B'$  do
5:      $B' \leftarrow B$ 
6:      $g \leftarrow \text{TrainClassifier}(B, g)$ 
7:      $B|S \leftarrow \text{Project}(B, g)$ 
8:     for block  $c$  in  $B|S$  do
9:       while  $B$  contains block  $b$  for which  $B \neq$ 
       $\text{Split}(b, c, B)$  do
10:         $B \leftarrow \text{Split}(b, c, B)$ 
11:      end while
12:    end for
13:  end while
14: end procedure

```

4.1 Abstraction

Algorithm 1 gives an overview of our abstraction process. Since we find MDP homomorphisms from experience, we first need to collect transitions that cover all regions of the state-action space. For simple environments, a random exploration policy provides such experience. But, a random walk is clearly not sufficient for more realistic environments because it rarely reaches the goal of the task. Therefore, we use the vanilla version of a deep Q-network [9] to collect the initial experience in bigger environments.

Subsequently, we partition the state-action space of the original MDP based on the collected experience with our Online Partition Iteration algorithm (Algorithm 2). The algorithm is described in detail in Subsection 4.2. The state-action partition B —the output of Algorithm 2—induces a quotient, or abstract, MDP according to Definition 2.6.

The quotient MDP enables both planning optimal actions for the current task (Subsection 4.2) and learning new tasks faster (Subsection 4.4).

4.2 Partitioning algorithm

Our online partitioning algorithm (Algorithm 2) is based on the Partition Iteration algorithm from [5]. It was originally developed for stochastic bisimulation based partitioning, and we adapted it to MDP homomorphisms (following Ravindran’s sketch [13]). Algorithm 4.2 starts with a reward respecting partition obtained by separating transitions that receive distinct rewards (*SplitRewards*). The reward respecting partition is subsequently refined with the *Split* (Algorithm 4) operation until a stopping condition is met. *Split*(b, c, B) splits a state-action block b from state-action partition B with respect to a state block c obtained by projecting the partition B onto the state space.

The projection of the state-action partition onto the state space (Algorithm 3) is the most complex component of our method. We train a classifier g , which can be an arbitrary model, to classify state-action pairs into their corresponding state-action blocks. The training set consists of all transitions the agent experienced, with

each transition belonging to a particular state-action block. During State Projection, g evaluates a state under a sampled set of actions, predicting a state-action block for each action. For discrete action spaces, the set should include all available actions. The set of predicted state-action blocks determines which state block the state belongs to.

Figure 2 illustrates the projection process: a single state s is evaluated under four actions: a_1, a_2, a_3 and a_4 . The first three actions are classified into the state-action block b_1 , whereas the last action is assigned to block b_3 . Therefore, s belongs to the state block identified by the set of the predicted state-action blocks $\{b_1, b_3\}$.

The output of Online Partition Iteration is a partition B of the state-action space Φ . According to Definition 2.6, the partition induces a quotient MDP. Since the quotient MDP is fully defined, we can compute its optimal Q-values with a dynamic programming method such as Value Iteration [17].

To be able to act according to the quotient MDP, we need to connect it to the original MDP in which we select actions. Given a current state s and a set of actions admissible in s, A_s , we predict the state-action block of each pair $(s, a_i), a_i \in A_s$ using the classifier g . Note that Online Partition Iteration trains g in the process of refining the partition. This process of predicting state-action block corresponds to a single step of State Projection: we determine which state block s belongs to. Since each state in the quotient MDP corresponds to a single state block (by Definition 2.6), we can map s to some state s' in the quotient MDP.

Given the current state s' in the quotient MDP, we select the action with the highest Q-value and map it back to the underlying MDP. An action in the quotient MDP can correspond to more than one action in the underlying MDP. For instance, an action that places a puck on the ground can be executed in many locations, while still having the same Q-value in the context of puck stacking. We break the ties between actions by sampling a single action in proportion to the confidence predicted by g : g predict a state-action block with some probability given a state-action pair.

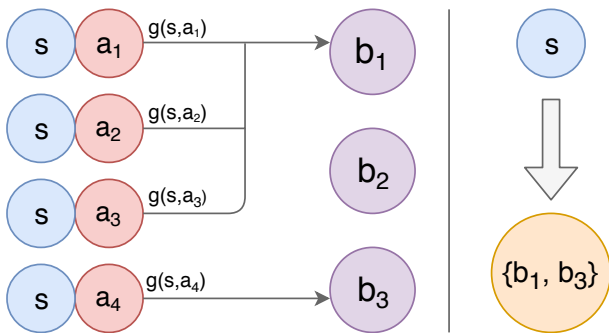


Figure 2: Projection (Algorithm 3) of a single state s . s is evaluated under actions a_1, a_2, a_3 and a_4 . For each pair (s, a_i) , the classifier g predicts its state-action block b_j . s belongs to a state block identified by the set of state-action blocks $\{b_1, b_3\}$.

Algorithm 3 State Projection

```

Input: State-action partition  $B$ , classifier  $g$ .
Output: State partition  $B|S$ .
1: procedure PROJECT
2:    $B|S \leftarrow \{\}$ 
3:   for block  $b$  in  $B$  do
4:     for transition  $t$  in  $b$  do
5:        $A_s \leftarrow \text{SampleActions}(t.\text{next\_state})$ 
6:        $B_s \leftarrow \{\}$ 
7:       for action  $a$  in  $A_s$  do
8:          $p \leftarrow g.\text{predict}(t.\text{next\_state}, a)$ 
9:          $B_s \leftarrow B_s \cup \{p\}$ 
10:      end for
11:      add  $t$  to  $B|S$  using  $B_s$  as the key
12:    end for
13:  end for
14: end procedure

```

4.3 Increasing robustness

Online Partition Iteration is sensitive to erroneous predictions by the classifier g . Since the collected transitions tend to be highly unbalanced and the mapping of state-action pairs into state-action blocks can be hard to determine, we include several augmentations that increase the robustness of our method. Some of them are specific to a neural network classifier.

- **class balancing:** The sets of state-action pairs belonging to different state-action blocks can be extremely unbalanced. Namely, the number of transitions that are assigned a positive reward is usually low. We follow the best practices from [2] and over-sample all minority classes so that the number of samples for each class is equal to the size of the majority class. We found decision trees do not require oversampling; hence, we use this method only with a neural network.
- **confidence calibration:** The latest improvements to neural networks, such as batch normalization [8] and skip connections [7] (both used by our neural network in Subsection 6.1), can cause miscalibration of the output class probabilities [6]. We calibrate the temperate of the softmax function applied to the output neurons using a multiclass version of Platt scaling [11] derived in [6]. The method requires a held-out validation set, which consists of 20% of all data in our case.
- **state-action block size threshold and confidence threshold:** During State Projection, the classifier g sometimes makes mistakes in classifying a state-action pair to a state-action block. Hence, the State Projection algorithm can assign a state to a wrong state block. This problems usually manifests itself with the algorithm "hallucinating" state blocks that do not exist in reality (note that there are $2^{\min\{|B|, |A|\}} - 1$ possible state blocks, given a state-action partition B). To prevent the *Split* function from over-segmenting the state-action partition due to these phantom state blocks, we only split a state-action block if the new blocks contain a number of samples higher than a threshold T_a . Furthermore, we exclude all predictions with confidence

Algorithm 4 Split

Input State-action block b , state block c , partition B .
Output State-action partition B' .

```

1: procedure SPLIT
2:    $b_1 \leftarrow \{\}, b_2 \leftarrow \{\}$ 
3:   for transition  $t$  in  $b$  do
4:     if  $transition.next\_state \in c$  then
5:        $b_1 \leftarrow b_1 \cup \{t\}$ 
6:     else
7:        $b_2 \leftarrow b_2 \cup \{t\}$ 
8:     end if
9:   end for
10:   $B' \leftarrow B$ 
11:  if  $|b_1| > 0 \ \&\& \ |b_2| > 0$  then
12:     $B' \leftarrow (B' \setminus \{b\}) \cup \{b_1, b_2\}$ 
13:  end if
14: end procedure

```

lower than some threshold T_c . Confidence calibration makes it easier to select the optimal value of T_c .

4.4 Transferring abstract MDPs

Solving a new task from scratch requires the agent to take a random walk before it stumbles upon a reward. The abstract MDP learned in the previous task can guide exploration by taking the agent into a starting state close to the goal of the task. However, how do we know which state block in the abstract MDP is a good start for solving a new task?

If we do not have any prior information about the structure of the next task, the agent needs to explore the starting states. To formalize this, we create $|B|S$ options, each taking the agent to a particular state in the quotient MDP from the first task. Each option is a tuple (I, π, β) with

- I being the set of all starting states of the MDP for the new task,
- π uses the quotient MDP from the previous task to select actions that lead to a particular state in the quotient MDP (see Subsection 4.2 for more details) and
- β terminates the option when the target state is reached.

The agent learns the Q -values of the options with a Monte Carlo update [17] with a fixed α (the learning rate)—the agent prefers options that make it reach the goal the fastest upon being executed. If the tasks are similar enough, the agent will find an option that brings it closer to the goal of the next task. If not, the agent can choose not to execute any option.

We use a deep Q-network to collect the initial experience in all transfer learning experiments. While our algorithm suffers from the same scalability issues as a deep Q-network when learning the *initial task*, our transfer learning method makes the learning of new tasks easier by guiding the agent's exploration.

5 PROOF OF CORRECTNESS

This section contains the proof of the correctness of our algorithm. We first prove two lemmas that support the main theorem. The

first lemma and corollary ensure that Algorithm 2 finds a reward respecting SSP partition.

LEMMA 5.1. *Given a reward respecting partition B of an MDP $M = \langle S, A, \Phi, P, R \rangle$ and $(s_1, a_1), (s_2, a_2) \in \Phi$ such that $T(s_1, a_1, [s']_{B|S}) \neq T(s_2, a_2, [s']_{B|S})$ for some $s' \in S$, (s_1, a_1) and (s_2, a_2) are not in the same block of any reward respecting SSP partition refining B .*

PROOF. Following the proof of Lemma 8.1 from [5]: proof by contradiction.

Let B' be a reward respecting SSP partition that is a refinement of B . Let $s' \in S$, $(s_1, a_1), (s_2, a_2) \in b \in B$ such that $T(s_1, a_1, [s']_{B|S}) \neq T(s_2, a_2, [s']_{B|S})$. Define B' such that $(s_1, a_1), (s_2, a_2)$ are in the same block and $[s']_{B|S} = \bigcup_{i=1}^k [s'_i]_{B'|S}$. Because B' is a reward respecting SSP partition, for each state block $[s'']_{B|S} \in B'|S$, $T(s_1, a_1, [s'']_{B|S}) = T(s_2, a_2, [s'']_{B|S})$. Then, $T(s_1, a_1, [s']_{B|S}) = \sum_{1 \leq i \leq k} T(s_1, a_1, [s'_i]_{B'|S}) = \sum_{1 \leq i \leq k} T(s_2, a_2, [s'_i]_{B'|S}) = T(s_2, a_2, [s']_{B|S})$. This contradicts $T(s_1, a_1, [s']_{B|S}) \neq T(s_2, a_2, [s']_{B|S})$. \square

COROLLARY 5.2. *Let B be a reward respecting partition of an MDP $M = \langle S, A, \Phi, P, R \rangle$, b a block in B and c a union of blocks from $B|S$. Every reward respecting SSP partition over Φ that refines B is a refinement of the partition $Split(b, c, B)$.*

PROOF. Following the proof of Corollary 8.2 from [5].

Let $c = \bigcup_{i=1}^n [s_i]_{B|S}$, $[s_i]_{B|S} \in B|S$. Let B' be a reward respecting SSP partition that refines B . $Split(b, c, B)$ will only split state-action pairs $(s_1, a_1), (s_2, a_2)$ if $T(s_1, a_1, c) \neq T(s_2, a_2, c)$. But if $T(s_1, a_1, c) \neq T(s_2, a_2, c)$, then there must be some k such that $T(s_1, a_1, c_k) \neq T(s_2, a_2, c_k)$ because for any $(s, a) \in \Phi$, $T(s, a, c) = \sum_{1 \leq m \leq n} T(s, a, c_m)$. Therefore, we can conclude by Lemma 5.1 that $[(s_1, a_1)]_{B'} \neq [(s_2, a_2)]_{B'}$. \square

The versions of Partition Iteration from [5] and [13] partition a fully-defined MDP. We designed our algorithm for the more realistic case, where only a stream of experience is available. This change makes the algorithm different only during State Projection (Algorithm 3). In the next lemma, we prove that the output of State Projection converges to a state partition as the number of experienced transitions goes to infinity.

LEMMA 5.3. *Let $M = \langle S, A, \Phi, P, R \rangle$ be an MDP with a finite A , a finite or infinite S , a state-action space Φ that is a separable metric space and a deterministic P defined such that each state-action pair is visited with a probability greater than zero. Let $SampleAction(s) = A_s, \forall s \in S$ (Algorithm 3, line 5). Let t_1, t_2, \dots be i.i.d. random variables that represent observed transitions, g a 1 nearest neighbor classifier that classifies state-action pairs into state-action blocks and let $(s, a)_n$ the nearest neighbor to (s, a) from a set of n transitions $X_n = \{t_1, t_2, \dots, t_n\}$. Let B_n be a state-action partition over X_n and $S_n = \bigcup_{t \in X_n} t.next_state$. Let $(B_n|S_n)'$ be a state partition obtained by the State Projection algorithm with g taking neighbors from X_n . $(B_n|S_n)' \rightarrow B_n|S_n$ as $n \rightarrow \infty$ with probability one.*

PROOF. $B_n|S_n$ is obtained by projecting B_n onto S_n . In this process, S_n is divided into blocks based on $B(s) = \{(s', a)]_{B_n} | (s', a) \in \Phi, s = s'\}$, the set of distinct blocks containing pairs of which s is a

Task	Options	Baseline	Baseline, share weights
2 puck stack to 3 puck stack	2558 ± 910	5335 ± 1540	10174 ± 5855
3 puck stack to 2 and 2 puck stack	2382 ± 432	-	3512 ± 518
2 puck stack to stairs from 3 pucks	2444 ± 487	4061 ± 1382	4958 ± 3514
3 puck stack to stairs from 3 pucks	1952 ± 606	4061 ± 1382	5303 ± 3609
2 puck stack to 3 puck component	2781 ± 605	3394 ± 999	6641 ± 5582
stairs from 3 pucks to 3 puck stack	3947 ± 873	5335 ± 1540	6563 ± 4299
stairs from 3 pucks to 2 and 2 puck stacks	5552 ± 3778	-	5008 ± 1998
stairs from 3 pucks to 3 puck component	3996 ± 2693	3394 ± 999	4856 ± 3600
3 puck component to 3 puck stack	3729 ± 742	5335 ± 1540	8540 ± 4908
3 puck component to stairs from 3 pucks	3310 ± 627	4061 ± 1382	2918 ± 328

Table 1: Transfer experiments in the pucks world domain. We measure the number of time steps before the agent reached the goal in at least 80% of episodes over a window of 50 episodes and report the mean and standard deviation over 10 trials. Unreported scores mean that the agent never reached this target. The column labeled *Options* represents our transfer learning method (Subsection 4.4), *Baseline* is deep Q-network described in Subsection 6.1 that does not retain any information from the initial task and *Baseline, share weights* copies the trained weights of the network from initial task to the transfer task. The bolded scores correspond to a statistically significant result for a Welch’s t-test with $P < 0.1$.

component, $s \in S_n$. Given $SampleAction(s) = A_s, \forall s \in S_n$, line 8 in Algorithm 3 predicts a $b \in B$ for each $(s', a) \in \Phi$, such that $s = s'$. By the Convergence of the Nearest Neighbor Lemma [3], $(s, a)_n$ converges to (s, a) with probability one. The rest of the Algorithm 3 exactly follows the projection procedure (Definition 3), therefore, $(B_n|S_n)' \rightarrow B_n|S_n$ with probability one. \square

Finally, we prove the correctness of our algorithm given an infinite stream of i.i.d. experience. While the i.i.d. assumption does not usually hold in reinforcement learning (RL), the deep RL literature often leverages the experience buffer [9] to ensure the training data is diverse enough. Our algorithm also contains a large experience buffer to collect the data needed to run Online Partition Iteration.

THEOREM 5.4 (CORRECTNESS). *Let $M = \langle S, A, \Phi, P, R \rangle$ be an MDP with a finite A , a finite or infinite S , a state-action space Φ that is a separable metric space and a deterministic P defined such that each state-action pair is visited with a probability greater than zero. Let $SampleAction(s) = A_s, \forall s \in S$ (Algorithm 3, line 5). Let t_1, t_2, \dots be i.i.d. random variables that represent observed transitions, g a 1 nearest neighbor classifier that classifies state-action pairs into state-action blocks. As the number of observed t_i goes to infinity, Algorithm 2 computes a reward respecting SSP partition over the observed state-action pairs with probability one.*

PROOF. Loosely following the proof of Theorem 8 from [5].

Let B be a partition over the observed state-action pairs, S the set of observed states and $(B|S)'$ the result of StateProjection(B, g) (Algorithm 3).

Algorithm 2 first splits the initial partition such that a block is created for each set of transitions with a distinct reward (line 2). Therefore, Algorithm 2 refines a reward respecting partition from line 2 onward.

Algorithm 1 terminates with B when $B = Split(b, [s]_{(B|S)'}, B)$ for all $b \in B, [s]_{(B|S)' \in (B|S)'$. $Split(b, [s]_{(B|S)'}, B)$ will split any block b containing $(s_1, a_1), (s_2, a_2)$ for which $T(s_1, a_1, [s]_{(B|S)'}) \neq T(s_2, a_2, [s]_{(B|S)'})$. According to Lemma 2, $(B|S)' \rightarrow B|S$ as $N \rightarrow \infty$ with probability one. Consequently, any partition returned by Algorithm 2 must be a reward respecting SSP partition.

Since Algorithm 2 first creates a reward respecting partition, and each step only refines the partition by applying $Split$, we can conclude by Corollary 1 that each partition encountered, including the resulting partition, must contain a reward respecting SSP partition. \square

6 EXPERIMENTS

We investigate the following questions with our experiments:

- (1) Can Online Partition Iteration find homomorphisms in environments with continuous state spaces and high-dimensional action spaces (characteristic for robotic manipulation tasks)?
- (2) Do options induced by quotient MDPs speed-up the learning of new tasks?
- (3) How does Online Partition Iteration compare to the only previous approach to finding homomorphisms[21]?

Section 6.1 describes our experiments and results concerning question 1 and 2, and section 6.2 presents a comparison with the prior work. We discuss the results in Subsection 6.3.

6.1 Continuous pucks world

We designed the pucks world domain (Figure 3) to approximate real-world robotic manipulation tasks. The state is represented by a 112x112 depth image and each pixel in the image is an admissible action. Hence, 12544 actions can be executed in each state. Environments with such a high branching factor favor homomorphisms, as they can automatically group actions into a handful of classes (e.g.

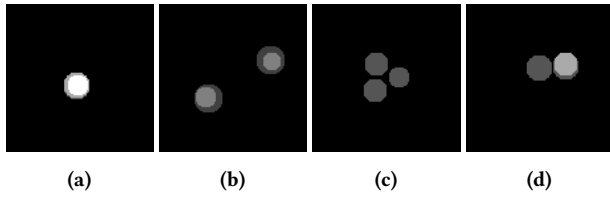


Figure 3: The goal states of the four types of tasks in our continuous pucks world domain. a) the task of stacking three pucks on top of one another, b) making two stacks of two pucks, c) arranging three pucks into a connected component, d) building stairs from three pucks.

"pick puck" and "do nothing") for each state. If an action corresponding to a pixel inside of a puck is selected, the puck is transported in the agent's hand. In the same way, the agent can stack pucks on top of each other or place them on the ground. Corner cases such as placing a puck outside of the environment or making a stack of pucks that would collapse are not allowed. The agent gets a reward of 0 for visiting each non-goal state and a reward of 10 for reaching the goal states. The environment terminates when the goal is reached or after 20 time steps. We implemented four distinct types of tasks: stacking pucks in a single location, making two stacks of pucks, arranging pucks into a connected component and building stairs from pucks. The goal states of the tasks are depicted in Figure 3. We can instantiate each task type with a different number of pucks, making the space of possible tasks and their combinations even bigger.

To gather the initial experience for partitioning, we use a shallow fully-convolutional version of the vanilla deep Q-network. Our implementation is based on the OpenAI baselines [4] with the standard techniques: separate target network with a weight update every 100 time steps and a replay buffer that holds the last 10 000 transitions. The network consists of five convolutional layers with the following settings (number of filters, filter sizes, strides): (32, 8, 4), (64, 8, 2), (64, 3, 1), (32, 1, 1), (2, 1, 1). The ReLU activation function is applied to the output of each layer except for the last one. The last layer predicts two maps of Q-values with the resolution 14x14 (for 112x112 inputs)—the two maps correspond to the two possible hand states: "hand full" and "hand empty". The appropriate map is selected based on the state of the hand, and bilinear upsampling is applied to get a 112x112 map of Q-values, one for each action. We trained the network with a Momentum optimizer with the learning rate set to 0.0001 and momentum to 0.9, batch size was set to 32. The agent interacted with the environment for 15000 episodes with an ϵ -greedy exploration policy; ϵ was linearly annealed from 1.0 to 0.1 for 40000 time steps.

Online Partition Iteration requires a second neural network—the classifier g . Our initial experiments showed that the predictions of the architecture described above lack in resolution. Therefore, we chose a deeper architecture: the DRN-C-26 version of Dilated Residual Networks [22]. We observed that the depth of the network is more important than the width (the number of filters in each layer) for our classification task. Capping the number of filters at 32 (the original architecture goes up to 512 filters in the last layers)

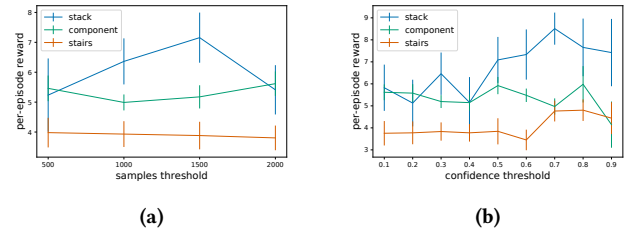


Figure 4: Grid searches over state-action block size thresholds and prediction confidence thresholds (described in Subsection 4.3). The y-axis represents the average per-episode reward (the maximum is 10) obtained by planning in the quotient MDP induced by the resulting partition. Stack, component and stairs represent the three puck world tasks shown in Figure 3. We report the means and standard deviations over 20 runs with different random seeds.

produces results indistinguishable from the original network. DRN-C-26 decreases the resolution of the feature maps in three places using strided convolutions, we downsample only twice to keep the resolution high. We train the network for 1500 steps during every iteration of Online Partition Iteration. The learning rate for the Momentum optimizer started at 0.1 and was divided by 10 at steps 500 and 1000, momentum was set to 0.9. The batch size was set to 64 and the weight decay to 0.0001.

Figure 4 reports the results of a grid search over state-action block size thresholds and classification confidence thresholds described in Subsection 4.3. Online Partition Iteration can create a near-optimal partition for the three pucks stacking task. On the other hand, our algorithm is less effective in the component arrangement and stairs building tasks. These two tasks are more challenging in terms of abstraction because the individual state-action blocks are not as clearly distinguishable as in puck stacking.

Next, we investigate if the options induced by the found partitions transfer to new tasks (Table 1). For puck stacking, a deep Q-network augmented with options from the previous tasks significantly outperforms both of our baselines. "Baseline" is a vanilla deep Q-network that does not retain any information from the initial task, whereas "Baseline, share weights" remembers the learned weights from the first task. Options are superior to the weights sharing baseline because they can take the agent to any desirable state, not just the goal. For instance, the 2 and 2 puck stacking task benefits from the option "make a stack of two pucks"; hence, options enable faster learning than weight sharing. We would also like to highlight one failure mode of the weight sharing baseline: the agent can sometimes get stuck repeatedly reaching the goal of the initial task without going any further. This behavior is exemplified in the transfer experiment from 2 puck stacking to 3 puck stacking. Here, the weight sharing agent continually places two pucks on top of one another, then lifts the top puck and places it back, which leads to slower learning than in the no-transfer baseline. Options do not suffer from this problem.

As reported in Figure 4, the learned partitions for the stairs building and component arrangement tasks underperform compared to puck stacking. Regardless, we observed a speed-up compared to the

no-transfer baseline in all experiments except for the transfer from stairs from 3 pucks to 3 puck component. Options also outperform weight sharing in 3 out of 5 experiments with the non-optimal partitions, albeit not significantly.

6.2 Discrete blocks world

Finally, we compare our partitioning algorithm to the decision tree method from [21] in the blocks world environment. The environment consists of three blocks that can be placed in four positions. The blocks can be stacked on top of one another, and the goal is to place a particular block, called the *focus block*, in a goal position and height. With four positions and three blocks, 12 tasks of increasing difficulty can be generated. The agent is penalized with -1 reward for each action that does not lead to the goal; reaching the goal state results in 100 reward.

Although a neural network can learn to solve this task, a decision tree trains two orders of magnitude faster and often reaches better performance. We used a decision tree from the scikit-learn package [10] with the default settings as our g classifier. All modifications from Subsection 4.3 specific to a neural network were omitted: class balancing and confidence thresholding. We also disabled the state-action block size threshold because the number of unique transitions generated by this environment is low and the decision tree does not make many mistakes. Despite the decision tree reaching high accuracy, we set a limit of 100 state-action blocks to avoid creating thousands of state-action pairs if the algorithm fails. The abstract MDP was recreated every 3000 time steps and the task terminated after 15000 time steps.

Figure 5 compares the decision tree version of our algorithm with the results reported in [21]. There are several differences between our experiments and the algorithm in [21]: Wolfe’s algorithm works with a Controlled Markov Process (CMP), an MDP augmented with an output function that provides richer supervision than the reward function. Therefore, their algorithm can start segmenting state-action blocks before it even observes the goal state. CMPs also allow an easy transfer of the learned partitions from one task to another; we solve each task separately. On the other hand, each action in Wolfe’s version of the task has a 0.2 chance of failure, but we omit this detail to satisfy the assumptions of our algorithm. Even though each version of the task is easier in some ways are harder in others, we believe the comparison with the only previous algorithm that solves the same problem is valuable.

6.3 Discussion

We return to the questions posed at the beginning of this section. Online Partition Iteration can find the right partition of the state-action space as long as the individual state-action blocks are clearly identifiable. For tasks with more complex classification boundaries, the partitions found are suboptimal, but still useful. We showed that options speed-up learning and outperform the baselines in the majority of the transfer experiments. Our algorithm also outperformed the only previous method of the same kind [21] in terms of finding a consistent partition.

The main drawback of Online Partition Iteration is that it is highly influenced by the accuracy of the classifier g . During state projection, it takes only one incorrectly classified action (out of

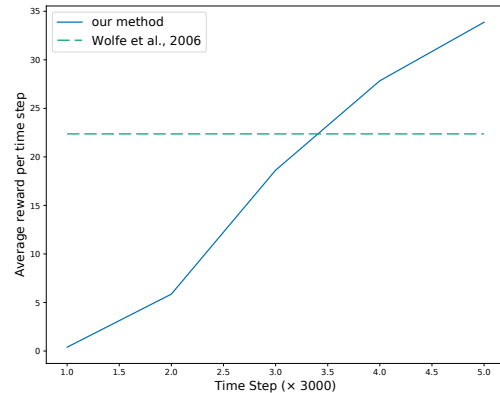


Figure 5: Comparison with Wolfe et al. [21] in the Blocks World environment. The horizontal line marks the highest mean reward per time step reached by Wolfe et al. We averaged our results over 100 runs with different goals.

12544 actions used in our pucks world experiments) for the state block classification to be erroneous. Confidence thresholding helps in the task of stacking pucks (Figure 4b), as it can filter out most of the errors. However, trained classifiers for the other two tasks, arranging components and building stairs, often produce incorrect predictions with a high confidence.

Moreover, the errors during state projection get amplified as the partitioning progresses. Note that the dataset of state-action pairs (inputs) and state-action blocks (classes) is created based on the previous state partition, which is predicted by the classifier. In other words, the version of the classifier g at step t generates the classes that will be used for its training at step $t + 1$. A classifier trained on noisy labels is bound to make even more errors at the next iteration. In particular, we observed that the error rate grows exponentially in the number of steps required to partition the state-action space.

In these cases, the partitioning algorithm often stops because of the limit on the number of state-action blocks (10 for the pucks domain). That is why the performance for the component arrangement and stairs building tasks is not sensitive to the state-action block size threshold (Figure 4a). Nevertheless, these noisy partitions also help with transfer learning, as shown in Table 1.

7 CONCLUSION

We developed Online Partition Iteration, an algorithm for finding abstract MDPs in discrete and continuous state spaces from experience, building on the existing work of Givan et al., Ravindran and Wolfe et al. [5, 13, 21]. We proved the correctness of our algorithm under certain assumptions and demonstrated that it could successfully abstract MDPs with high-dimensional continuous state spaces and high-dimensional discrete action spaces. In addition to being interpretable, the abstract MDPs can guide exploration when learning new tasks. We created a transfer learning method in the framework of options [18], and demonstrated that it outperforms the baselines in the majority of experiments.

REFERENCES

- [1] Richard Bellman. 1957. A Markovian Decision Process. *Journal of Mathematics and Mechanics* 6, 5 (1957), 679–684.
- [2] Mateusz Buda, Atsuto Maki, and Maciej A. Mazurowski. 2018. A systematic study of the class imbalance problem in convolutional neural networks. *Neural networks: the official journal of the International Neural Network Society* 106 (2018), 249–259.
- [3] T. Cover and P. Hart. 1967. Nearest neighbor pattern classification. *IEEE Transactions on Information Theory* 13, 1 (January 1967), 21–27. <https://doi.org/10.1109/TIT.1967.1053964>
- [4] Prafulla Dhariwal, Christopher Hesse, Oleg Klimov, Alex Nichol, Matthias Plappert, Alec Radford, John Schulman, Szymon Sidor, Yuhuai Wu, and Peter Zhokhov. 2017. OpenAI Baselines. <https://github.com/openai/baselines>. (2017).
- [5] Robert Givan, Thomas Dean, and Matthew Greig. 2003. Equivalence notions and model minimization in Markov decision processes. *Artificial Intelligence* 147, 1 (2003), 163 – 223. [https://doi.org/10.1016/S0004-3702\(02\)00376-4](https://doi.org/10.1016/S0004-3702(02)00376-4) Planning with Uncertainty and Incomplete Information.
- [6] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. 2017. On Calibration of Modern Neural Networks. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*. 1321–1330. <http://proceedings.mlr.press/v70/guo17a.html>
- [7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep Residual Learning for Image Recognition. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016), 770–778.
- [8] Sergey Ioffe and Christian Szegedy. 2015. Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. In *Proceedings of the 32nd International Conference on Machine Learning (Proceedings of Machine Learning Research)*, Francis Bach and David Blei (Eds.), Vol. 37. PMLR, Lille, France, 448–456. <http://proceedings.mlr.press/v37/ioffe15.html>
- [9] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharmarajan Kumar, Daan Wierstra, Shane Legg, and Demis Hassabis. 2015. Human-level control through deep reinforcement learning. *Nature* 518 (25 Feb 2015), 529 EP –.
- [10] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Courville, M. Brucher, M. Perrot, and E. Duchesnay. 2011. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research* 12 (2011), 2825–2830.
- [11] John C. Platt. 1999. Probabilistic Outputs for Support Vector Machines and Comparisons to Regularized Likelihood Methods. In *ADVANCES IN LARGE MARGIN CLASSIFIERS*. MIT Press, 61–74.
- [12] S. Rajendran and M. Huber. 2009. Learning to generalize and reuse skills using approximate partial policy homomorphisms. In *2009 IEEE International Conference on Systems, Man and Cybernetics*. 2239–2244. <https://doi.org/10.1109/ICSMC.2009.5345891>
- [13] Balaraman Ravindran. 2004. *An Algebraic Approach to Abstraction in Reinforcement Learning*. Ph.D. Dissertation. AAI3118325.
- [14] Balaraman Ravindran and Andrew G Barto. 2004. Approximate homomorphisms: A framework for non-exact minimization in Markov decision processes. (2004).
- [15] Vishal Soni and Satinder Singh. 2006. Using Homomorphisms to Transfer Options Across Continuous Reinforcement Learning Domains. In *Proceedings of the 21st National Conference on Artificial Intelligence - Volume 1 (AAAI'06)*. AAAI Press, 494–499.
- [16] Jonathan Sorg and Satinder Singh. 2009. Transfer via Soft Homomorphisms. In *Proceedings of The 8th International Conference on Autonomous Agents and Multiagent Systems - Volume 2 (AAMAS '09)*. International Foundation for Autonomous Agents and Multiagent Systems, Richland, SC, 741–748.
- [17] Richard S. Sutton and Andrew G. Barto. 1998. *Introduction to Reinforcement Learning* (1st ed.). MIT Press, Cambridge, MA, USA.
- [18] Richard S. Sutton, Doina Precup, and Satinder Singh. 1999. Between MDPs and semi-MDPs: A framework for temporal abstraction in reinforcement learning. *Artificial Intelligence* 112, 1 (1999), 181 – 211. [https://doi.org/10.1016/S0004-3702\(99\)00052-1](https://doi.org/10.1016/S0004-3702(99)00052-1)
- [19] Jonathan J. Taylor, Doina Precup, and Prakash Panangaden. 2008. Bounding Performance Loss in Approximate MDP Homomorphisms. In *Proceedings of the 21st International Conference on Neural Information Processing Systems (NIPS'08)*. Curran Associates Inc., USA, 1649–1656.
- [20] Alicia P. Wolfe. 2006. Defining Object Types and Options Using MDP Homomorphisms.
- [21] Alicia Peregrin Wolfe and Andrew G. Barto. 2006. Decision Tree Methods for Finding Reusable MDP Homomorphisms. In *Proceedings of the 21st National Conference on Artificial Intelligence - Volume 1 (AAAI'06)*. AAAI Press, 530–535.
- [22] Fisher Yu, Vladlen Koltun, and Thomas A. Funkhouser. 2017. Dilated Residual Networks. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2017), 636–644.