

Affordance-Aware Planning

Abstract

Planning algorithms for non-deterministic domains are often intractable in large state spaces due to the well-known curse of dimensionality. Existing approaches to planning in large stochastic state spaces fail to prevent autonomous agents from considering many actions that are obviously irrelevant to a human solving the same task. To prevent agents from applying irrelevant actions we formalize the notion of *affordances* as state space independent, goal-oriented knowledge added to an Object Oriented Markov Decision Process (OO-MDP). Affordances prune irrelevant actions based on the agent's goal and the current state, reducing the number of state-action pairs the planner must evaluate in order to formulate a near optimal policy. Moreover, affordances may be provided by an expert or may be learned without supervision. We demonstrate our approach in the state-rich Minecraft domain, showing significant increases in speed, reductions in state space exploration, and improvements in the quality of the synthesized policy. Additionally, we show that learned affordances often surpass the performance of those provided by experts. We also demonstrate that affordance-aware planning enables a Baxter robot to assist a person performing a cooking task.

Introduction

Robots operating in unstructured, stochastic environments such as a factory floor or a kitchen face a difficult planning problem due to the large state space and inherent uncertainty due to unreliable perception and actuation [5, 16]. Robotic planning tasks are often formalized as a stochastic sequential decision making problem, modeled as a Markov Decision Process (MDP) [28]. In these problems, the agent must find a mapping from states to actions for some subset of the state space that enables the agent to achieve a goal while minimizing costs along the way. However, many robotics tasks are so complex that modeling them as an MDP results in a massive state-action space, which in turn restricts the types of robotics problems that are computationally tractable. For example, when a robot is manipulating objects in an environment, an object can be placed anywhere in a large set of locations. The size of the state space increases exponentially with the number of objects, which bounds the placement problems that the robot is able to expediently solve.

Moreover, the difficulty of the task is compounded by the fact that most of these objects and locations are irrelevant. For instance, when making brownies, the oven and flour are important, while the soy sauce and sauté pan are not. For a different task, such as stir-frying broccoli, a different set of objects and actions are relevant.

To address this state-action space explosion, prior work has explored adding knowledge to the planner, such as options [27] and macroactions [6, 22]. However, while these methods allow the agent to search more deeply in the state space, they add non-primitive actions to the planner which *increase* the size of the state-action space. The resulting augmented space is even larger, which can have the paradoxical effect of increasing the search time for a good policy [15]. Deterministic forward-search algorithms like hierarchical task networks (HTNs) [21], and temporal logical planning (TLPlan) [2, 3], add knowledge to the planner that greatly increases planning speed, but do not generalize to stochastic domains. Additionally, the knowledge provided to the planner by these methods is quite extensive, reducing the agent's autonomy.

Instead, we augment an Object Oriented Markov Decision Process (OO-MDP) with a formalization of *affordances*. Affordances were originally proposed by Gibson [12] as action possibilities prescribed by an agent's capabilities in an environment. We rigorously formalize the notion of an affordance as knowledge added to an OO-MDP that prunes irrelevant actions on a state by state basis. Our affordances can be specified by hand and also learned through experience, making them a concise, transferable, and learnable means of representing useful planning knowledge. Our experiments demonstrate that affordances provide dramatic improvements for a variety of planning tasks compared to baselines, and apply across different state spaces. Moreover, while manually provided affordances outperform baselines, affordances learned through experience yield even greater improvements. We conduct experiments in the game Minecraft, which has a very large state-action space, and on a real-world robotic cooking assistant.

Technical Approach

An MDP is a five-tuple: $\langle \mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \gamma \rangle$, where \mathcal{S} is a state space; \mathcal{A} is the agent's set of actions; \mathcal{T} denotes $\mathcal{T}(s' | s, a)$, the transition probability of an agent applying action $a \in \mathcal{A}$

in state $s \in \mathcal{S}$ and arriving in $s' \in \mathcal{S}$; $\mathcal{R}(s, a, s')$ denotes the reward received by the agent for applying action a in state s and transitioning to state s' ; and $\gamma \in [0, 1)$ is a discount factor that defines how much the agent prefers immediate rewards over future rewards (the agent prefers to maximize immediate rewards as γ decreases).

Our representation of affordances builds on Object-Oriented MDPs (OO-MDP) [10]. An OO-MDP efficiently represents the state of an MDP through the use of objects and predicates. An OO-MDP state is a collection of objects, $O = \{o_1, \dots, o_o\}$. Each object o_i belongs to a class, $c_j \in \{c_1, \dots, c_c\}$. Every class has a set of attributes, $Att(c) = \{c.a_1, \dots, c.a_a\}$, each of which has a domain, $Dom(c.a)$, of possible values. OO-MDPs enable planners to use predicates over classes of objects. That is, the OO-MDP definition also includes a set of predicates \mathcal{P} that operate on the state of objects to provide additional high-level information about the MDP state.

OO-MDP predicates provide state space independence. For a given planning domain, OO-MDP objects often appear across tasks. Since predicates operate on collections of objects, they generalize beyond specific state spaces within the given domain. For instance, in Minecraft, a predicate checking the contents of the agent’s inventory generalizes beyond any particular Minecraft task. We capitalize on this state space independence by using OO-MDP predicates as features for action pruning.

Modeling the Optimal Actions

Our goal is to formalize affordances in a way that enables a planning algorithm to prune away suboptimal actions in each state. We define the optimal action set, \mathcal{A}^* , for a given state s and goal G as:

$$\mathcal{A}^* = \{a \mid Q_G^*(s, a) = V_G^*(s)\}, \quad (1)$$

where, $Q_G^*(s, a)$ and $V_G^*(s)$ represent the optimal Q function and value function, respectively.

We aim to learn a probability distribution over the optimality of each action for a given state (s), goal (G), and knowledge base (K) from which action pruning may be informed. Thus, we want to infer a Bernouli for each action’s optimality:

$$\Pr(a_i \in \mathcal{A}^* \mid s, G, K) \quad (2)$$

for $i \in \{1, \dots, |\mathcal{A}|\}$, where \mathcal{A} is the OO-MDP action space.

We formalize our knowledge base, K , as a set of n paired preconditions and goal types, $\{(p_1, g_1) \dots (p_n, g_n)\}$, along with a parameter vector, θ . We abbreviate each pair (p_j, g_j) to δ_j for simplicity. Each precondition $p \in \mathcal{P}$ is a *predicate* in predicate space, \mathcal{P} , defined by the OO-MDP, and $g \in \mathcal{G}$ is a *goal type* in goal space. For example, a predicate might be *nearTrench(agent)* which is true when the agent is standing near a trench. We assume that the goal space consists of logical expressions of state predicates. A goal type specifies the sort of problem the agent is trying to achieve. In the context of Minecraft, a goal type might refer to the agent retrieving an object of a certain type from the environment, reaching a particular location, or crafting an object or structure.

We rewrite Equation 2 replacing K with its constituents:

$$\begin{aligned} \Pr(a_i \in \mathcal{A}^* \mid s, G, K) \\ = \Pr(a_i \in \mathcal{A}^* \mid s, G, \delta_1 \dots \delta_n, \theta_i) \end{aligned} \quad (3)$$

where θ_i represents the set of parameters relevant to modeling the probability of action $a_i \in \mathcal{A}^*$.

We introduce the indicator function f , which returns 1 if and only if δ_j ’s predicate is true in the provided state s , and δ_j ’s goal type matches the agent’s current goal, G :

$$f(\delta, s, G) = \begin{cases} 1 & \delta.p(s) \wedge (G == \delta.g) \\ 0 & \text{otherwise} \end{cases} \quad (4)$$

Evaluating f for each δ_j given the current state and goal gives rise to a set of binary features, $\phi_j = f(\delta_j, s, G)$, which we use to reformulate our probability distribution.

$$\begin{aligned} \Pr(a_i \in \mathcal{A}^* \mid s, G, \delta_1 \dots \delta_n, \theta_i) \\ = \Pr(a_i \in \mathcal{A}^* \mid \phi_1, \dots, \phi_n, \theta_i) \end{aligned} \quad (5)$$

This distribution may be modeled in a number of ways making this approach quite flexible. However, to enable efficient learning, we model our distribution using Naive Bayes. First we factor using Bayes’ rule:

$$= \frac{\Pr(\phi_1 \dots \phi_n, \mid a_i \in \mathcal{A}^*, \theta_i) \Pr(a_i \in \mathcal{A}^* \mid \theta_i)}{\Pr(\phi_1, \dots, \phi_n \mid \theta_i)} \quad (6)$$

Next we assume that each feature is conditionally independent of the others, given whether the action is optimal:

$$= \frac{\prod_{j=1}^n \Pr(\phi_j \mid a_i \in \mathcal{A}^*, \theta_i) \Pr(a_i \in \mathcal{A}^* \mid \theta_i)}{\Pr(\phi_1, \dots, \phi_n \mid \theta_i)} \quad (7)$$

Finally, we define the prior on the optimality of each action to be the fraction of the time each action was optimal during training. This distribution fully describes the model used by our affordance-aware planner.

Learning the Optimal Actions

Our approach to modeling the optimality of each action allows affordances to be learned through unsupervised experience. To learn affordances, we provide a set of training worlds (W) for which the optimal policy, π , may be tractably computed. Then, we compute the maximum likelihood estimate of the parameter vector θ_i for each action.

Under our Bernouli Naive Bayes model, we estimate the parameters $\theta_{i,0} = \Pr(a_i)$ and $\theta_{i,j} = \Pr(\phi_j \mid a_i)$, for $j \in \{1, \dots, n\}$, where the maximum likelihood estimates are:

$$\theta_{i,0} = \frac{C(a_i)}{C(a_i) + C(\bar{a}_i)} \quad (8)$$

$$\theta_{i,j} = \frac{C(\phi_j, a_i)}{C(a_i)} \quad (9)$$

Here, $C(a_i)$ is the number of observed occurrences where a_i was optimal across all worlds W , $C(\bar{a}_i)$ is the number of observed occurrences where a_i was not optimal, and $C(\phi_j, a_i)$ is the number of occurrences where $\phi_j = 1$ and

a_i was optimal. In all cases, optimality was determined according to π .

$$C(a_i) = \sum_{w \in W} \sum_{s \in w} (\pi(s) == a_i) \quad (10)$$

$$C(\bar{a}_i) = \sum_{w \in W} \sum_{s \in w} (\pi(s) \neq a_i) \quad (11)$$

$$C(\phi_j, a_i) = \sum_{w \in W} \left(\sum_{s \in w} \pi(s) == a_i \wedge \phi_j == 1 \right) \quad (12)$$

Affordance-Aware Planning

We complete our formalization of affordances by exploring three different ways to prune actions based on the distribution on the optimality of each action.

First, affordances may be specified by a domain expert in place of the Naive Bayes. In this approach, an expert specifies a set of actions associated with a precondition-goal type pair. When the precondition is active and goal type is the same as the agent’s present goal, the actions suggested by the affordance are included in the agent’s action set. For instance, if an agent is standing above a block of buried gold and is trying to smelt a block of gold, then an expert may indicate that the agent should consider the actions of looking down and digging. Table 1 shows several examples of expert defined affordances. All actions contributed by active affordances are grouped to yield the set of actions to consider for each state.

Second, actions may be pruned by thresholding the posterior. In this method, the affordances remove any actions whose probability of being optimal is below the provided threshold for each state. The threshold was determined empirically, and was set to $\frac{0.2}{|A|}$, where $|A|$ is the size of the full action space of the OO-MDP. This threshold is quite conservative, and means that our approach only prunes actions which are extremely unlikely to be optimal.

Lastly, actions may be pruned by sampling actions from the probability distribution as specified by Equation 7. We treat each action’s probability mass as a Bernoulli trial and sample across the entire action set. In preliminary results, this method did not perform as well as baselines - likely because the weights associated with each action were too small. In future work, we are interested in investigating more sophisticated approaches to pruning actions with sampling.

Through the use of any of the above methods, an affordance-aware planner prunes actions on a state by state basis, focusing the agent on relevant action possibilities of the environment, consequently reducing planning time. Any planner operating in an OO-MDP may be made affordance-aware with this approach.

In a recent review on the theory of affordances, Chemero [7] suggests that an affordance is a relation between the features of an environment and an agent’s abilities. Our approach grounds this interpretation, where the features of the environment correspond to the goal-dependent state features, ϕ , and the agent’s abilities correspond to the OO-MDP action set. In our model, there is an affordance for each δ_j ,



Figure 1: A gold smelting task in the Minecraft domain. The agent’s goal is to mine a block of gold, move to the forge and then smelt the gold in the forge to produce gold ingots. **stefie10: Can you label the goal and the forge in the picture?**

with preconditions $\delta_j.p$, goal type $\delta_j.g$ and action distribution $\Pr(a_i \in \mathcal{A}^* | \phi_j, \theta)$, which is computed in our Naive Bayes model by marginalizing over all the features not associated with ϕ_j .

Results

We evaluate our approach using the game Minecraft and a collaborative robotic cooking task. Minecraft is a 3-D blocks game in which the user can place, craft, and destroy blocks of different types. Minecraft’s physics and action space allow users to create complex systems, including logic gates and functional scientific graphing calculators¹. Minecraft serves as a model for robotic tasks such as cooking assistance, assembling items in a factory, object retrieval, and complex terrain traversal. As in these tasks, the agent operates in a very large state-action space in an uncertain environment. Figure 1 shows a scene from one of our Minecraft problems. **label: Stephen: describe your domain here**

Minecraft Tests

Our experiments consisted of five common tasks in Minecraft, including constructing bridges over trenches, smelting gold, tunneling through walls, basic path planning, and digging to find an object. We tested on randomized worlds of varying size and difficulty. The generated test worlds varied in size from tens of thousands of states to hundreds of thousands of states.

For learned affordances, a knowledge base was derived from the training data, which consisted of 20 simple state

¹<https://www.youtube.com/watch?v=wgJfVRhotlQ>

Precondition	Goal Type	Actions
lookingTowardGoal	atLocation	{move}
lavaInFront	atLocation	{rotate}
lookingAtGold	hasGoldOre	{destroy}

Table 1: Examples of Expert Provided Affordances

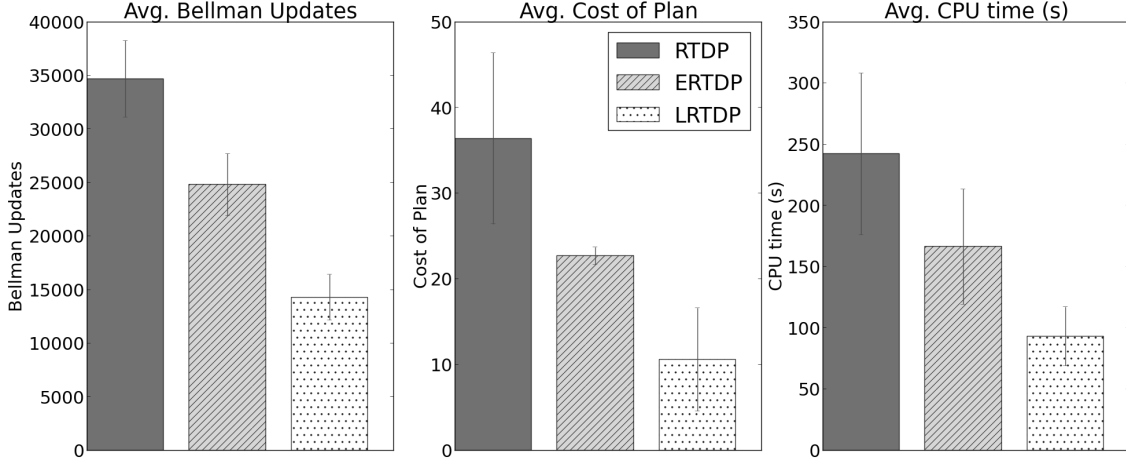


Figure 2: Average Results From All Maps

spaces of each map type (100 total maps), each approximately a 1,000-10,000 state world. We conducted all tests with a single knowledge base.

We use Real Time Dynamic Programming (RTDP) [4] as our baseline planner, a sampling-based algorithm that does not require the planner to visit all states. We compare RTDP with learned affordance-aware RTDP (LRTDP), and expert-defined affordance-aware RTDP (ERTDP). LRTDP pruned actions according to the thresholded pruning method. We terminated each planner when the maximum change in the value function was less than 0.01 for 100 consecutive policy rollouts, or the planner failed to converge after 1000 rollouts. We set the reward function to -1 for all transitions, except transitions to states in which the agent was in lava, where we set the reward to -10 . The goal was set to be terminal. The discount factor was set to $\lambda = 0.99$. For all experiments, movement actions (move, rotate, jump) had a small probability (0.05) of incorrectly applying a different movement action.

The evaluation metrics for each trial were the number of Bellman updates executed by each planning algorithm, the accumulated reward of the average plan, and the CPU time taken to find a plan. Table 3 shows the average Bellman updates, accumulated reward, and CPU time for RTDP, LRTDP and ERTDP after planning in 20 different maps of each goal type (100 total).

Planner	Bellman	Reward	CPU
RTDP	27439 (± 2348)	-22.6 (± 9)	107 (± 33)
LRTDP	9935 (± 1031)	-12.4 (± 1)	53 (± 5)
RTDP+Opt	26663 (± 2298)	-17.4 (± 4)	129 (± 35)
LRTDP+Opt	9675 (± 953)	-11.5 (± 1)	93 (± 10)
RTDP+MA	31083 (± 2468)	-21.7 (± 5)	336 (± 28)
LRTDP+MA	9854 (± 1034)	-11.7 (± 1)	162 (± 17)

Table 2: Affordances vs. Temporally Extended Actions

Figure 2 shows the results averaged across all maps. The full results for these trials are shown in Table 3. Because the planners were forced to terminate after only 1000 rollouts, they did not always converge to the optimal policy. As the results suggest, LRTDP on average found a comparably better plan (10.6 cost) than ERTDP (22.7 cost) and RTDP (36.4 cost), found the plan in significantly fewer Bellman updates (14287.5 to ERTDP’s 24804.1 and RTDP’s 34694.3) and in less CPU time (93.1s to ERTDP’s 166.4s and RTDP’s 242.0s). These results indicate that while learned affordances gave the largest improvements, expert-provided affordances can also significantly enhance performance, and, depending on the domain, could add significant value in making large state-spaces tractable without the overhead of supplying training worlds.

Temporally Extended Actions and Affordances

We compared our approach to Temporally Extended Actions: macroactions and options. We conducted these experiments with the same configurations as our Minecraft experiments. Domain experts provided the option policies and macroactions.

Table 2 indicates the results of comparing RTDP equipped with macro actions, options, and affordances across 100 dif-

ferent executions in the same randomly generated Minecraft worlds. The results are averaged across tasks of each type presented in Table 3. As the results suggest, both macroactions and options add a significant amount of time to planning. This increase is because it is computationally expensive to predict the expected reward associated with applying an option or a macroaction. Furthermore, the branching factor of the state-action space significantly increases when augmented with additional actions, causing the planner to run for longer and perform more Bellman updates. With affordances, the planner found a better plan in less CPU time, and with fewer Bellman updates. These results support the claim that affordances can handle the augmented action space provided by temporally extended actions by pruning away unnecessary actions.

ARTDP and Baxter

label: Insert description of baxter stuff here



Figure 3: Placeholder for baxter results/image

Related Work

In this section, we discuss the differences between affordance-aware planning and other forms of knowledge engineering that have been used to accelerate planning. We divide these approaches into those that are built to plan in stochastic domains, and those that are designed for use with deterministic domains.

Stochastic Approaches

Here, we compare other approaches of action pruning and knowledge engineering that provide speedups to planners in stochastic domains.

Temporally Extended Actions Temporally extended actions are actions that the agent can select like any other action of the domain, except executing them results in multiple primitive actions being executed in succession. Two common forms of temporally extended actions are *macro-actions* [14] and *options* [27]. Macro-actions are actions that always execute the same sequence of primitive actions. Options are defined with high-level policies that accomplish specific sub tasks. For instance, when an agent is near a door, the agent can engage the ‘door-opening-option-policy’, which switches from the standard high-level planner to running a policy that is hand crafted to open doors. Although the classic options framework is not generalizable to different state spaces, creating *portable* options is a topic of active research [19, 17, 24, 8, 1, 18].

Since temporally extended actions may negatively impact planning time [15] by adding to the number of actions the

Planner	Bellman	Reward	CPU
<i>Mining Task</i>			
RTDP	17142.1 (± 3843)	-6.5 (± 1)	17.6s (± 4)
ERTDP	14357.4 (± 3275)	-6.5 (± 1)	31.9s (± 8)
LRTDP	12664.0 (± 9340)	-12.7 (± 5)	33.1s (± 23)
<i>Smelting Task</i>			
RTDP	30995.0 (± 6730)	-8.6 (± 1)	45.1s (± 14)
ERTDP	28544.0 (± 5909)	-8.6 (± 1)	72.6s (± 19)
LRTDP	2821.9 (± 662)	-9.8 (± 2)	7.5s (± 2)
<i>Wall Traversal Task</i>			
RTDP	45041.7 (± 11816)	-56.0 (± 51)	68.7s (± 22)
ERTDP	32552.0 (± 10794)	-34.5 (± 25)	96.5s (± 39)
LRTDP	24020.8 (± 9239)	-15.8 (± 5)	80.5s (± 34)
<i>Trench Traversal Task</i>			
RTDP	16183.5 (± 4509)	-8.1 (± 2)	53.1s (± 22)
ERTDP	8674.8 (± 2700)	-8.2 (± 2)	35.9s (± 15)
LRTDP	11758.4 (± 2815)	-8.7 (± 1)	57.9s (± 20)
<i>Plane Traversal Task</i>			
RTDP	52407 (± 18432)	-82.6 (± 42)	877.0s (± 381)
ERTDP	32928 (± 14997)	-44.9 (± 34)	505.3s (± 304)
LRTDP	19090 (± 9158)	-7.8 (± 1)	246s (± 159)

Table 3: RTDP vs. Affordance-Aware RTDP

agent can choose from in a given state, combining affordances with temporally extended actions allows for even further speedups in planning, as demonstrated in Table 2. In other words, affordances are complementary knowledge to options and macroactions.

Action Pruning Sherstov and Stone [26] considered MDPs with a very large action set and for which the action set of the optimal policy of a source task could be transferred to a new, but similar, target task to reduce the learning time required to find the optimal policy in the target task. The main difference between our affordance-based action set pruning and this action transfer work is that affordances prune away actions on a state by state basis, whereas the learned action pruning is on per task level. Further, with goal types, affordances may be attached to subgoal planning for a significant benefit in planning tasks where complete subgoal knowledge is known.

Rosman and Ramamoorthy [25] provide a method for learning action priors over a set of related tasks. Specifically, they compute a Dirichlet distribution over actions by extracting the frequency that each action was optimal in each state for each previously solved task.

Action priors can only be used with planning/learning algorithms that work well with an ϵ -greedy rollout policy, while affordances can be applied to almost any planning algorithm. In addition, action priors are only utilized for

fraction ϵ of the time steps, which is typically quite small, limiting the improvement they can make to the planning speed. Finally, as variance in tasks explored increases, the priors will become more uniform. In contrast, affordance-aware planning can handle a wide variety of tasks in a single knowledge base, as demonstrated by Table 3.

Heuristics Heuristics in MDPs are used to convey information about the value of a given state-action pair with respect to the task being solved and typically take the form of either *value function initialization*, or *reward shaping*. Initializing the value function to an admissible close approximation of the optimal value function has been shown to be effective for LAO* and RTDP [13]. Reward shaping is an alternative approach to providing heuristics. The planning algorithm uses a modified version of the reward function that returns larger rewards for state-action pairs that are expected to be useful, but does not guarantee convergence to an optimal policy unless certain properties of the shaped reward are satisfied [23].

However, heuristics are highly dependent on the reward function and state space of the task being solved, whereas affordances are state space independent and may be learned easily for different reward functions. If a heuristic can be provided, the combination of heuristics and affordances may even more greatly accelerate planning algorithms than either approach alone.

Deterministic Approaches

There have been several successful attempts at engineering knowledge to decrease planning time for deterministic planners. These are fundamentally solving a different problem from what we are interested in since they deal with non-stochastic problems, but there are a number of salient parallels and contrasts to be drawn nonetheless.

Hierarchical Task Networks Traditional Hierarchical Task Networks (HTNs) employ *task decompositions* to aid in planning [11]. The goal at hand is decomposed into smaller tasks which are in turn decomposed into smaller tasks. This decomposition continues until primitive tasks that are immediately achievable are derived. The current state of the task decomposition, in turn, informs constraints which reduce the space over which the planner searches. At a high level both HTNs and affordances fulfill the same role: both achieve action pruning by exploiting some form of supplied knowledge.

However there are two essential distinctions between affordances and traditional HTNs. (1) HTNs do not incorporate reward into their planning. Consequently, they lack mathematical guarantees of optimal planning. (2) On a qualitative level, the degree of supplied knowledge in HTNs surpasses that of affordances: whereas affordances simply require relevant propositional functions, HTNs require not only constraints for sub-tasks but a hierarchical framework of arbitrary complexity.

Temporal Logic Bacchus and Kabanza [2, 3] provided planners with domain dependent knowledge in the form of a first-order version of linear temporal logic (LTL), which they

used for control of a forward-chaining planner. With this methodology, a STRIPS style planner may be guided through the search space by checking whether candidate plans do not falsify a given knowledge base of LTL formulas, often achieving polynomial time planning in exponential space.

ellis: Maybe add how LTL is similar to what we are doing?

The primary difference between this body of work and affordance-aware planning is that affordances may be learned increasing autonomy of the agent and flexibility of the approach, while LTL formulas are far too complicated to learn effectively, placing dependence on an expert.

Conclusion

label: Conclusion could use some work/rewriting We proposed a novel approach to represent transferable knowledge in terms of *affordances* [12]. Our affordances allow an agent to efficiently prune actions based on learned knowledge, providing a significant reduction in the number of state-action pairs the agent needs to evaluate in order to act near optimally. We demonstrated the effectiveness of the affordance model by comparing RTDP to its affordance-aware equivalents in a series of challenging planning tasks in the Minecraft domain. Further, we designed a full learning process that allows an agent to autonomously learn useful affordances that may be used across a variety of task types, reward functions, and state spaces, allowing for convenient extensions to robotic applications.

Additionally, we compared the effectiveness of augmenting planners with affordances, with temporally extended actions, and both. The results suggest that affordances may be combined with temporally extended actions to provide improvements in planning, since affordances can handle the extended action space.

Lastly, we deployed an affordance-aware planner on a robot in a collaborative cooking task with a massive state space. **label: More baxter details here.**

In the future, we hope to automatically discover useful state space specific subgoals online - a topic of some active research [20, 9]. Since affordances are goal specific, they will plug easily into high-level subgoal planning. Automatic discovery of subgoals would allow affordance-aware planners to take advantage of this connection, and would further reduce the size of the explored state-action space by improving the effectiveness of action pruning.

Additionally, we hope to explore additional methods that capitalize on the distribution over optimal actions, such as incorporating affordances with a forward search sparse sampling algorithm [29], or replacing the Naive Bayes model with a more sophisticated model, such as Logistic Regression or a Noisy-OR. We are also investigating methods of learning the thresholded value in a more principled way - one such approach is to initialize the planner with a strict threshold, and slowly relax the threshold until a near optimal policy is found. We are also interested in updating model parameters on-line by using planning data to update the distribution over optimal actions.

References

- [1] D. Andre and S.J. Russell. State abstraction for programmable reinforcement learning agents. In *Eighteenth national conference on Artificial intelligence*, pages 119–125. American Association for Artificial Intelligence, 2002.
- [2] Fahiem Bacchus and Froduald Kabanza. Using temporal logic to control search in a forward chaining planner. In *In Proceedings of the 3rd European Workshop on Planning*, pages 141–153. Press, 1995.
- [3] Fahiem Bacchus and Froduald Kabanza. Using temporal logics to express search control knowledge for planning. *Artificial Intelligence*, 116:2000, 1999.
- [4] Andrew G Barto, Steven J Bradtko, and Satinder P Singh. Learning to act using real-time dynamic programming. *Artificial Intelligence*, 72(1):81–138, 1995.
- [5] Mario Bollini, Stefanie Tellex, Tyler Thompson, Nicholas Roy, and Daniela Rus. Interpreting and executing recipes with a cooking robot. In *Proceedings of International Symposium on Experimental Robotics (ISER)*, 2012.
- [6] Adi Botea, Markus Enzenberger, Martin Müller, and Jonathan Schaeffer. Macro-ff: Improving ai planning with automatically learned macro-operators. *Journal of Artificial Intelligence Research*, 24:581–621, 2005.
- [7] Anthony Chemero. An outline of a theory of affordances. *Ecological psychology*, 15(2):181–195, 2003.
- [8] T. Croonenborghs, K. Driessens, and M. Bruynooghe. Learning relational options for inductive transfer in relational reinforcement learning. *Inductive Logic Programming*, pages 88–97, 2008.
- [9] Özgür Şimşek, Alicia P. Wolfe, and Andrew G. Barto. Identifying useful subgoals in reinforcement learning by local graph partitioning. In *Proceedings of the 22Nd International Conference on Machine Learning, ICML '05*, pages 816–823, New York, NY, USA, 2005. ACM. ISBN 1-59593-180-5. doi: 10.1145/1102351.1102454. URL <http://doi.acm.org/10.1145/1102351.1102454>.
- [10] C. Diuk, A. Cohen, and M.L. Littman. An object-oriented representation for efficient reinforcement learning. In *Proceedings of the 25th international conference on Machine learning, ICML '08*, 2008.
- [11] Kutluhan Erol, James Hendler, and Dana S Nau. Htn planning: Complexity and expressivity. In *AAAI*, volume 94, pages 1123–1128, 1994.
- [12] JJ Gibson. The concept of affordances. *Perceiving, acting, and knowing*, pages 67–82, 1977.
- [13] Eric A Hansen and Shlomo Zilberstein. Solving markov decision problems using heuristic search. In *Proceedings of AAAI Spring Symposium on Search Techniques from Problem Solving under Uncertainty and Incomplete Information*, 1999.
- [14] Milos Hauskrecht, Nicolas Meuleau, Leslie Pack Kaelbling, Thomas Dean, and Craig Boutilier. Hierarchical solution of markov decision processes using macro-actions. In *Proceedings of the Fourteenth conference on Uncertainty in artificial intelligence*, pages 220–229. Morgan Kaufmann Publishers Inc., 1998.
- [15] Nicholas K. Jong. The utility of temporal abstraction in reinforcement learning. In *Proceedings of the Seventh International Joint Conference on Autonomous Agents and Multiagent Systems*, 2008.
- [16] Ross A. Knepper, Stefanie Tellex, Adrian Li, Nicholas Roy, and Daniela Rus. Single assembly robot in search of human partner: Versatile grounded language generation. In *Proceedings of the HRI 2013 Workshop on Collaborative Manipulation*, 2013.
- [17] G. Konidaris and A. Barto. Efficient skill learning using abstraction selection. In *Proceedings of the Twenty First International Joint Conference on Artificial Intelligence*, pages 1107–1112, 2009.
- [18] G. Konidaris, I. Scheidwasser, and A. Barto. Transfer in reinforcement learning via shared features. *The Journal of Machine Learning Research*, 98888:1333–1371, 2012.
- [19] George Konidaris and Andrew Barto. Building portable options: Skill transfer in reinforcement learning. In *Proceedings of the International Joint Conference on Artificial Intelligence, IJCAI '07*, pages 895–900, January 2007.
- [20] Amy McGovern and Andrew G. Barto. Automatic discovery of subgoals in reinforcement learning using diverse density. In *In Proceedings of the eighteenth international conference on machine learning*, pages 361–368. Morgan Kaufmann, 2001.
- [21] Dana Nau, Yue Cao, Amnon Lotem, and Hector Munoz-Avila. Shop: Simple hierarchical ordered planner. In *Proceedings of the 16th International Joint Conference on Artificial Intelligence - Volume 2, IJCAI'99*, pages 968–973, San Francisco, CA, USA, 1999. Morgan Kaufmann Publishers Inc. URL <http://dl.acm.org/citation.cfm?id=1624312.1624357>.
- [22] M Newton, John Levine, and Maria Fox. Genetically evolved macro-actions in ai planning problems. *Proceedings of the 24th UK Planning and Scheduling SIG*, pages 163–172, 2005.
- [23] Andrew Y Ng, Daishi Harada, and Stuart Russell. Policy invariance under reward transformations: Theory and application to reward shaping. In *ICML*, volume 99, pages 278–287, 1999.
- [24] Balaraman Ravindran and Andrew Barto. An algebraic approach to abstraction in reinforcement learning. In *Twelfth Yale Workshop on Adaptive and Learning Systems*, pages 109–144, 2003.
- [25] Benjamin Rosman and Subramanian Ramamoorthy. What good are actions? accelerating learning using learned action priors. In *Development and Learning and Epigenetic Robotics (ICDL), 2012 IEEE International Conference on*, pages 1–6. IEEE, 2012.
- [26] A.A. Sherstov and P. Stone. Improving action selection in mdp's via knowledge transfer. In *Proceedings of the 20th national conference on Artificial Intelligence*, pages 1024–1029. AAAI Press, 2005.
- [27] Richard S Sutton, Doina Precup, and Satinder Singh. Between mdps and semi-mdps: A framework for temporal abstraction in reinforcement learning. *Artificial intelligence*, 112(1):181–211, 1999.
- [28] Sebastian Thrun, Wolfram Burgard, and Dieter Fox. *Probabilistic robotics*. MIT Press, 2008.
- [29] Thomas Walsh, Sergiu Goschin, and Michael Littman. Integrating sample-based planning and model-based reinforcement learning, 2010. URL <http://www.aaai.org/ocs/index.php/AAAI/AAAI10/paper/view/1880>.