
Knowledge Representation and Interactive Learning of Domain Knowledge for Human-Robot Interaction

Mohan Sridharan

M.SRIDHARAN@BHAM.AC.UK

School of Computer Science, University of Birmingham, Birmingham, UK B15 2TT

Ben Meadows

BMEA011@AUCKLANDUNI.AC.NZ

Department of Electrical and Computer Engineering, The University of Auckland, Auckland, NZ

Abstract

This paper describes an integrated architecture for representing, reasoning with, and interactively learning domain knowledge in the context of human-robot collaboration. Answer Set Prolog, a non-monotonic logical reasoning paradigm, is used to represent and reason with incomplete commonsense domain knowledge, computing a plan for any given goal and diagnosing unexpected observations. ASP-based reasoning is also used to guide the interactive learning of previously unknown actions as well as axioms that encode affordances, action preconditions, and effects. This learning takes as input observations from active exploration, reactive action execution, and human (verbal) descriptions, and the learned actions and axioms are used for subsequent reasoning. The architecture is evaluated on a simulated robot assisting humans in an indoor domain.

1. Introduction

Consider robots¹ in an office or warehouse delivering objects to particular locations and people. Information about such domains often includes commonsense knowledge, especially default knowledge that holds in all but a few exceptional circumstances, e.g., “books are usually in the library, but cookbooks are in the kitchen”. Domain knowledge may also include some information about action preconditions and effects, and action capabilities, i.e., affordances. In complex domains, human participants may lack the time and expertise to provide comprehensive domain knowledge or elaborate feedback, so robots will need to reason with incomplete domain knowledge and revise this knowledge over time. The architecture described in this paper is a step towards these capabilities. It implements the following theoretical tenets:

- Knowledge elements symbolically encode objects; relations modeling domain attributes and actions at different levels of abstraction; and axioms composed of these relations.
- Knowledge elements are revised non-monotonically by reasoning with knowledge and observed outcomes of actions that may be immediate or delayed.

1. We use the terms “robot”, “agent”, and “learner” interchangeably in this paper.

- Affordances are relations defined jointly over the attributes of agents and the attributes of objects in the context of particular actions.
- Reasoning, learning, and interaction are coupled; values of state-action pairs are revised using observations obtained from active exploration and reactive action execution.

The combination of these tenets is novel, and our implementation exploits the complementary strengths of declarative programming, probabilistic reasoning, and interactive learning. We explored subsets of these tenets in prior work (Sridharan & Meadows, 2017; Sridharan et al., 2017b). Here, we focus on the interplay between representation, reasoning and learning, and abstract away some aspects of our architecture: we merge some levels of representation and do not probabilistically model perceptual uncertainty. We describe the following capabilities of the architecture:

- Incomplete domain knowledge described in an action language is translated into a relational representation in Answer Set Prolog (ASP) for inference, planning and diagnostics. ASP-based reasoning also limits interactive learning to the relevant part of the domain.
- Previously unknown actions’ names, preconditions, effects, and objects over which they operate, along with the associated affordances, are learned using decision-tree induction and relational reinforcement learning by processing observations from active exploration, reactive action execution, and verbal cues from humans.

The novelty in comparison with our prior work lies in learning actions along with the preconditions, effects and affordances, and in doing so interactively through relational inference and reinforcement based on active exploration, reactive action execution, and human (verbal) inputs. We evaluate these capabilities on a simulated robot delivering objects to particular people or places in an indoor domain. For instance, a robot equipped with our architecture and assisting humans in an office could learn interactively that it cannot pick up heavy objects with its electromagnetic arm, and that grasping a brittle cup will cause the cup to break. We measure the robot’s ability to acquire missing knowledge, and the resultant ability to compute minimal and correct plans for assigned goals. We begin by reviewing related work (Section 2) and describing our architecture (Section 3). We present an experimental evaluation of the framework (Section 4) and offer closing thoughts (Section 5).

2. Related Work

Agents deployed in complex domains often have to represent and reason with incomplete domain knowledge, and learn from observations. Early work used a first-order logic representation and incrementally refined the action operators but did not allow for different outcomes in different contexts (Gil, 1994). Such approaches have difficulties with non-monotonic logical reasoning and merging new, unreliable information with existing beliefs. Research in logics has provided non-monotonic reasoning formalisms, e.g., with ASP (Erdem & Patoglu, 2012; Erdem et al., 2016). Researchers have combined ASP with inductive learning to monotonically learn causal laws (Otero, 2003), and expanded the theory of actions in ASP to learn and revise system descriptions or domain knowledge (Balduccini, 2007; Law et al., 2018). Architectures such as SOAR reason with hierarchical knowledge in first-order logic and process perceptual information probabilistically to acquire domain knowledge (Laird, 2012). Cognitive architectures based on non-monotonic logic, or

a combination of logic and probabilistic representations, have also been used to support reasoning and learning in robotics (Scheutz et al., 2007; Zhang et al., 2015). However, approaches based on classical first-order logic are often not expressive enough. Algorithms based on logic programming, on the other hand, find it difficult to support one or more of the desired capabilities, such as efficient and incremental learning of knowledge, learning from interactions, and reasoning with large probabilistic components. Existing algorithms and architectures also do not support generalization over learned knowledge as described in this paper.

Many formalizations have been proposed for representing, reasoning with, and learning affordances (Zech et al., 2017). Existing approaches represent affordances as possible effects of actions or behaviors (Guerin et al., 2013), or as emergent, functional, and/or contextual properties based on attributes of the domain and the objects (Sarathy & Scheutz, 2016). These approaches have used logics, probabilistic reasoning or a combination of both. We recently developed existing work to provide a theory of affordances for cognitive systems (Langley et al., 2018), making claims related to (a) representing knowledge of affordances; (b) using this knowledge for plan understanding or plan generation; and (c) acquiring the knowledge of affordances. The architecture described in this paper explores some of these claims, e.g., that affordances are relations defined jointly over attributes of objects and agents in the context of specific actions.

Interactive task learning is a framework for acquiring domain knowledge using labeled examples or reinforcement signals obtained from domain observations, demonstrations, or human instructions (Chai et al., 2018; Laird et al., 2017). It can be viewed as building on early work on joint search through the space of hypotheses and observations (Simon & Lea, 1974). Algorithms for interactive task learning can be used to learn action models or search control rules that prevent particular actions from being considered under certain circumstances. Dynamic domains often require a learner to explore actions whose outcomes may be delayed, and Reinforcement learning (RL) elegantly addresses the credit assignment problem in such cases. RL is also appropriate for learning different types of knowledge that direct attention during both planning and learning. Approaches such as relational reinforcement learning (RRL) have been developed for efficient generalization in complex domains (Driessens & Ramon, 2003; Tadepalli et al., 2004). However, interactive learning algorithms, including those that exploit relational representations and different function approximations, tend to focus on a single planning task at a time and do not support the commonsense reasoning capabilities desired in robotics (Bloch & Laird, 2017; Driessens & Ramon, 2003). One exception was our prior work that combined the principles of non-monotonic logical reasoning and RRL to discover domain axioms (Sridharan & Meadows, 2017; Sridharan et al., 2017b). The architecture described in this paper significantly extends these capabilities to support interactive learning of actions along with their preconditions, effects, and affordances, based on relational inference and reinforcement using active exploration, reactive action execution, and verbal inputs.

3. The Reasoning Architecture

The architecture that forms the basis of our work uses commonsense domain knowledge, including a description of domain dynamics in an action language, to construct tightly-coupled representations at two resolutions, with the fine-resolution representation defined as a refinement of the coarse-resolution representation (Sridharan et al., 2017a). For any given goal, coarse-resolution reasoning

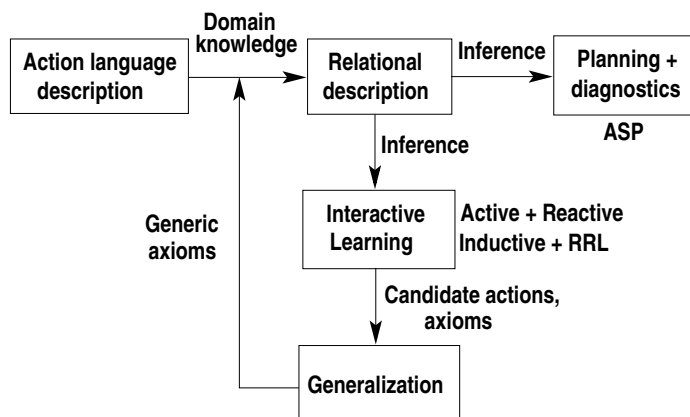


Figure 1. Non-monotonic logical inference with domain knowledge is used for (a) planning and diagnostics; and (b) guiding interactive learning of generic versions of affordances, actions, and related axioms based on observations obtained through active exploration and reactive action execution.

with commonsense knowledge provides a plan of abstract actions. Each abstract action is implemented as a sequence of concrete actions using a partially observable Markov decision process that reasons probabilistically over the relevant part of the fine-resolution representation. The observations obtained during this implementation are added to the coarse-resolution history for subsequent reasoning. In this paper, we abstract away the reasoning at different resolutions and the probabilistic modeling of perceptual uncertainty, to better focus on the interplay between representation, reasoning, and learning. As shown in Figure 1, the system translates the relational representation into an ASP program for planning and diagnostics. ASP-based reasoning also guides the interactive learning of actions, affordances, and the preconditions and effects of actions, using observations from active exploration, reactive execution, and human (verbal) input. The system uses the learned knowledge for subsequent reasoning. We illustrate these capabilities with the following domain.

Example 1 [*Robot Assistant (RA) Domain*] A simulated robot finds, labels, and delivers objects to people or places. Each place may have instances of objects such as *desk*, *book*, *cup* and *computer*. Each human has a *role* (e.g., *engineer*, *manager*, *salesperson*). Object attributes include *weight* (*heavy*, *light*), *surface* (*brittle*, *hard*), *status* (*intact*, *damaged*), and *labeled* (*true*, *false*). The robot’s arm has a *type* (*electromagnetic*, *pneumatic*). The robot’s actions include *pickup*, *putdown*, *move*, *label*, and *serve*. Some actions or axioms governing domain dynamics may be unknown, for example:

- A pneumatic arm cannot be used to serve a brittle object.
- Serving an object to a salesperson causes it to be labeled.
- Object with a brittle surface cannot be labeled unless the robot has an electromagnetic arm.

For simplicity, any other robots in the domain are assumed to have identical capabilities and cannot communicate with the learner. Humans and the learner can observe these robots. Humans can verbally describe other robots’ activities, e.g., “Robot labeled the hard, hefty item” to help the

learner acquire previously unknown knowledge. This domain becomes complex as the number of ground instances of objects and their attributes increases, e.g., there were $\approx 18,000$ combinations of ground static attributes and ≈ 11 million combinations of ground fluent terms in an instantiation of the domain that we use in our experiments.

3.1 Knowledge Representation and Reasoning

We first describe the action language encoding of the domain dynamics, and its translation to CR-Prolog programs for knowledge representation and reasoning.

Action Language. Action languages are formal models of parts of natural language used for describing transition diagrams of dynamic systems. Our architecture uses action language \mathcal{AL}_d (Gelfond & Inlezan, 2013) to describe the transition diagrams of the domain. \mathcal{AL}_d has a sorted signature with *statics*, i.e., domain attributes whose truth values cannot be changed by actions, *fluents*, i.e., domain attributes whose truth values can be changed by actions, and *actions*, a set of elementary operations. Fluents can be *basic*, which obey inertia laws and can be changed by actions, or *defined*, which do not obey the laws of inertia and are not changed directly by actions. A domain attribute or its negation is a *literal*. \mathcal{AL}_d allows three types of statements:

$$\begin{aligned} a \text{ causes } l_b \text{ if } p_0, \dots, p_m & \quad (\text{Causal law}) \\ l \text{ if } p_0, \dots, p_m & \quad (\text{State constraint}) \\ \text{impossible } a_0, \dots, a_k \text{ if } p_0, \dots, p_m & \quad (\text{Executability condition}) \end{aligned}$$

where a is an action, l is a literal, l_b is a basic literal, and p_0, \dots, p_m are domain literals.

Domain Representation: Signature and Axioms. The domain representation consists of system description \mathcal{D} , which is a collection of statements of \mathcal{AL}_d , and history \mathcal{H} . \mathcal{D} has a sorted signature Σ and axioms that describe the transition diagram τ . Σ defines the basic sorts, and domain attributes and actions. For the RA domain, basic sorts such as *place*, *robot*, *object*, *entity*, *book*, *weight*, and *step* (for temporal reasoning) are arranged hierarchically, e.g., *robot* and *object* are subsorts of *entity*. Σ also includes ground instances of sorts, e.g., *office*, *workshop*, *kitchen*, and *library* are instances of sort *place*. Domain attributes and actions are described in terms of the sorts of their arguments. The fluents include *loc(entity, place)*, the location of the robot and objects, with the locations of humans and other robots (if any) modeled as defined fluents whose values are obtained from external sensors; and *in_hand(robot, object)*, which denotes whether a particular object is in the robot’s hand. Static attributes include *arm_type(robot, type)*, *obj_status(object, status)* etc. Actions include *move(robot, place)*, *pickup(robot, object)*, and *serve(robot, object, person)*. The signature Σ also includes a relation *holds(fluent, step)* that implies a particular fluent is true at a particular timestep.

Axioms of the RA domain capture the domain’s dynamics. These axioms include causal laws, state constraints and executability conditions. For example:

$$\begin{aligned} \text{move}(\text{rob}_1, L) \text{ causes } \text{loc}(\text{rob}_1, L) \\ \text{serve}(\text{rob}_1, O, P) \text{ causes } \text{in_hand}(P, O) \end{aligned}$$

$$loc(O, L) \text{ if } loc(rob_1, L), in_hand(rob_1, O)$$

$$\text{impossible } pickup(rob_1, O) \text{ if } loc(rob_1, L_1), loc(O, L_2), L_1 \neq L_2$$

The history \mathcal{H} of a dynamic domain is usually a record of fluents observed to be true or false at a particular time step, i.e., $obs(fluent, boolean, step)$, and the occurrence of an action at a particular time step, i.e., $occurs(action, step)$. This notion was expanded to represent defaults describing the values of fluents in the initial state (Sridharan et al., 2017a), e.g., “books are usually in the library and if not there, they are normally in the office” is encoded as:

$$\text{initial default } loc(X, library) \text{ if } book(X)$$

$$\text{initial default } loc(X, office) \text{ if } book(X), \neg loc(X, library)$$

We can also encode exceptions to these defaults, e.g., “cookbooks are in the kitchen”.

Domain Representation: Affordances. We define affordances, i.e., action capabilities, as relations between attributes of robot(s) and object(s) in the context of particular actions. Negative (i.e., forbidding or dis-) affordances describe unsuitable combinations of objects, robots, and actions. Positive affordances describe permissible uses of objects in actions by agents, including exceptions to the executability conditions. We represent affordances in a distributed manner, as follows:

$$\text{impossible } A \text{ if } aff_forbids(ID, A)$$

$$aff_forbids(id_i, A) \text{ if } \dots$$

$$\text{impossible } A \text{ if } \dots, not\ aff_permits(ID, A)$$

$$aff_permits(id_j, A) \text{ if } \dots$$

The first two statements imply that action A cannot occur if it is not afforded, and specify the conditions (i.e., attributes of robot and objects) under which the action is not afforded. The last two statements imply that action A that is not considered during planning due to an executability condition may have a positive affordance as an exception, and encode the exception. Note that each action can have multiple indexed affordances; as discussed later, such a distributed representation improves generalization and simplifies inference.

Reasoning with Domain Knowledge. The reasoning tasks of the robot associated with a system description \mathcal{D} and history \mathcal{H} include planning, diagnostics and inference. To perform these tasks, the domain representation is translated into a program $\Pi(\mathcal{D}, \mathcal{H})$ in CR-Prolog², a variant of ASP that incorporates consistency restoring (CR) rules (Balduccini & Gelfond, 2003). ASP is based on stable model semantics, and supports *default negation* and *epistemic disjunction*, e.g., unlike “ $\neg a$ ” that implies *a is believed to be false*, “*not a*” only implies *a is not believed to be true*, and unlike “ $p \vee \neg p$ ” in propositional logic, “*p or not p*” is not tautologous. ASP can represent recursive definitions and constructs that are difficult to express in classical logic formalisms, and it supports non-monotonic logical reasoning, i.e., the ability to revise previously held conclusions based on new evidence. The program Π includes the signature and axioms of \mathcal{D} , inertia axioms, reality checks, and

2. We use the terms “ASP” and “CR-Prolog” interchangeably in this paper.

observations, actions, and defaults from \mathcal{H} . Every default also has a CR rule that allows the robot to assume the default’s conclusion is false under exceptional circumstances, to restore consistency. Π also includes other appropriate definitions and axioms, e.g., for planning, we add a definition of goal, and axioms that force the system to identify and include actions in the plan until the goal is achieved. Algorithms for computing the entailment, and for planning and diagnostics, then reduce these tasks to computing *answer sets* of the CR-Prolog programs. Each such answer set represents the set of inferred beliefs of an agent associated with the program.

Reasoning with incomplete or incorrect knowledge may overlook valid plans, find suboptimal plans, or provide plans whose execution has unintended outcomes. For instance, suppose the robot in the RA domain is asked to deliver textbook $book_1$ to the *office*. It uses default knowledge to compute a plan comprising four actions: (i) $move(rob_1, library)$; (ii) $pickup(rob_1, book_1)$; (iii) $move(rob_1, office)$; and (iv) $putdown(rob_1, book_1)$. However, this plan does not succeed because, unknown to the robot, its electromagnetic arm cannot pick up the heavy book. This exemplifies the kind of knowledge we seek to enable the robot to learn. We next describe the interactive learning of such knowledge.

3.2 Interactive Learning

In complex domains, learning previously unknown actions and related axioms in their most generic forms may require many labeled examples, and it may be difficult to obtain such labeled examples. Also, humans may not have the time and expertise required to provide labeled examples, and an action’s effects may be observed immediately or after a delay. To address these challenges, our architecture includes two schemes for interactive acquisition of labeled examples:

- (i) active learning from verbal cues provided by humans; and
- (ii) relational reinforcement learning that considers delayed rewards to mimic cumulative learning based on observations from active exploration or reactive action execution.

Learning from Human Interaction. The first interactive learning scheme uses verbal input provided by humans describing the observed behaviors of other robots in the domain. This scheme makes the following assumptions:

- Other robots in the domain have the same sensing and actuation capabilities as the learner;
- The human description of observed behavior focuses on one action at a time; also, this description may be ambiguous but it is not intentionally incorrect; and
- The learner can generate logic statements corresponding to the observed attributes of robot(s) or object(s) in the domain.

These assumptions are reasonable for many robotics domains, and they significantly simplify the learner’s interaction with humans.

The learner solicits human input when available and receives a transcribed verbal description of an action and observations of the action’s consequences, e.g., the learner may receive “The robot is labeling the fairly big textbook.” and $labeled(book_1)$. We use the Stanford log-linear part-of-speech (POS) tagger (Toutanova et al., 2003). We employ a left, second-order sequence information model to determine each word’s POS tag and append it to the word. In our example, the

output is a string such as “The_DT robot_NN is_VBZ labeling_VBG the_DT fairly_RB big_JJ text-book_NN”, where “VB” represents a verb, “NN” is a noun etc. The learner transforms this string to $\langle \text{word}, \text{POS} \rangle$ pairs, and transforms the sentence’s verb into first-person present-tense using rules from a lemma list (Someya, 1998) and WordNet (Miller, 1995), e.g., $\langle \text{is}, \text{VBZ} \rangle \langle \text{labeling}, \text{VBG} \rangle$ becomes the verb “label”. The learner also marks each noun phrase as a sequence of zero or more adjectival terms followed by a noun, discarding other interleaved words. Our example sentence’s noun phrases are *robot* and *big textbook*. Nouns signify object sorts and adjectival terms signify values of static attributes. To determine terms’ referents, WordNet relations such as *linked synsets* are used to find a synonym that is also a domain symbol, e.g., “big” and “heavy” share a WordNet synset, *heavy* is an attribute value, and $\text{book}(\text{book}_1)$ and $\text{obj_weight}(\text{book}_1, \text{heavy})$ are domain attributes. The matched sorts and attributes refer to particular objects. This approach requires static attributes’ values to be disjoint sets, and each noun phrase to signify an existing object.

Next, the robot constructs an action literal $\text{label}(\text{rob}_1, \text{book}_1)$ from the verb and object referents, and the arguments’ lowest-level sorts are assumed to indicate the action’s signature, e.g., $\text{label}(\#robot, \#book)$. If this candidate action does not match any known action literal, the robot lifts the literal, its arguments and the observed action consequences. This forms the basis for constructing candidate causal laws and generalizing over time. For instance:

$$\text{label}(\text{rob}_1, \text{book}_1) \text{ causes } \text{labeled}(\text{book}_1)$$

is lifted to:

$$\text{label}(R, B) \text{ causes } \text{labeled}(B)$$

On the other hand, if the new literal matches an existing one, the the common ancestor of each argument’s sort is considered, e.g., if the learner finds that newly discovered $\text{label}(\#robot, \#cup)$ matches with $\text{label}(\#robot, \#book)$, it will generalize to $\text{label}(\#robot, \#object)$. This learning scheme adapts existing natural language processing methods to work with our representation. It helps the learner acquire a previously unknown action’s signature and a related causal law. However, the learner may still not know axioms that govern the domain dynamics related to this action. This missing knowledge is acquired as described below.

Relational Reinforcement Learning. The second learning scheme enables axiom discovery by active exploration of particular transitions, or by exploration in response to unexpected and unexplained transitions.

Basic RL Formulation: To explore a particular transition, the resultant state is set as the goal of a reinforcement learning (RL) problem, i.e., the objective is to find state-action pairs most likely to lead to this (and other similar) states. The underlying MDP is defined by a set of states (S), set of actions (A), state transition function $T_f : S \times A \times S' \rightarrow [0, 1]$, and reward function $R_f : S \times A \times S' \rightarrow \mathfrak{R}$. Similar to classical RL formulations, T_f and R_f are unknown to the agent. Each state has ground atoms formed of the domain attributes (i.e., fluent terms and statics), and a boolean literal describing whether the most recent action had the expected outcome. Each action is a ground action of the system description. In robotics domains, the functions T_f and R_f are typically constructed from statistics collected in an initial training phase; T_f is a probabilistic model of the uncertainty in state

transitions, while R_f provides instantaneous rewards for executing particular actions in particular states. The RL formulation is constructed automatically from the system description—Sridharan et al. (2017a) describe a method for translating an ASP-based system description to a representation for probabilistic sequential decision making.

The values of state-action pairs are estimated in a series of episodes, until convergence, using the Q-learning algorithm (Sutton & Barto, 1998). In each episode, the agent executes a sequence of actions chosen by an ϵ -greedy algorithm with eligibility traces. The combinations of states and actions invalidated by existing axioms are not explored. An episode terminates when a time limit is exceeded or the target action succeeds. The physical configuration of objects is then reset to its state from the beginning of the episode, and a new episode begins.

Reducing Search Space: The basic RL formulation becomes computationally intractable for complex domains, even with more sophisticated RL algorithms. We address this problem in our architecture by automatically restricting learning to the object constants, domain attributes and axioms *relevant* to the transition under consideration. In other words, only parts of the system description that influence or are influenced by the transition of interest are included in the MDP underlying the RL formulation, which significantly reduces the search space. This notion of relevance is based on the following desiderata regarding the relations that may appear in a discovered axiom:

- For any static attribute that may exist in the body of the discovered axiom, we wish to explore all possible elements in the range of the attribute, e.g., for action $serve(rob_1, cup_1, person_1)$, all possible weights of cup_1 and roles of $person_1$ are explored.
- For any fluent that may appear in the body of the axiom, we wish to explore only those elements in the range of the fluent that occur in the state before or after the state transition. Any other element cannot, by design, be influenced by this transition anyway.

ASP-based reasoning is used to implement these broad requirements and automatically construct the system description $\mathcal{D}(T)$, the part of \mathcal{D} relevant to the transition T . To do so, we first define the object constants relevant to the transition of interest. These definitions are adapted from the definitions introduced in Sridharan et al. (2017a).

Definition 1 [*Relevant object constants*]

Let $T = \langle \sigma_1, a_{tg}, \sigma_2 \rangle$ be the transition of interest. Let $relCon(T)$ be the set of object constants of Σ of \mathcal{D} identified using the following rules:

1. Object constants from a_{tg} are in $relCon(T)$;
2. If $f(x_1, \dots, x_n, y)$ is a literal formed of a domain attribute, and the literal belongs to σ_1 or σ_2 , but not both, then x_1, \dots, x_n, y are in $relCon(T)$;
3. If the body B of an axiom of a_{tg} contains an occurrence of $f(x_1, \dots, x_n, Y)$, a term whose domain is ground, and $f(x_1, \dots, x_n, y) \in \sigma_1$, then x_1, \dots, x_n, y are in $relCon(T)$.

Object constants from $relCon(T)$ are said to be *relevant* to T . Consider σ_1 with $loc(rob_1, office)$, $loc(cup_1, office)$, and $loc(person_1, office)$, and $a_{tg} = serve(rob_1, cup_1, person_1)$. The object constants relevant to T include $rob_1, cup_1, person_1$, and $office$.

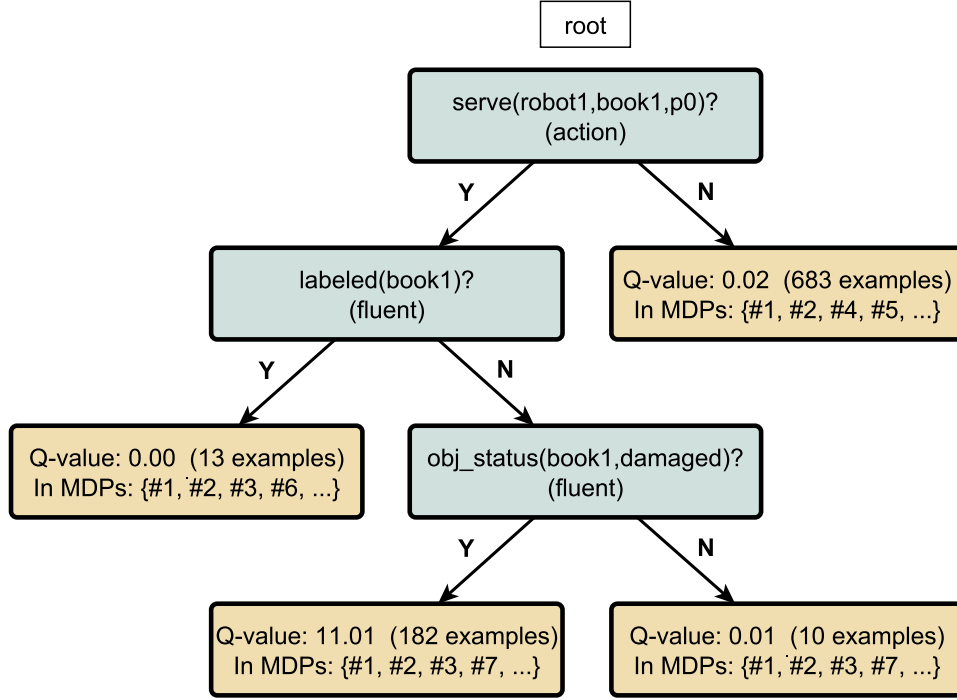


Figure 2. Illustrative example of a binary decision tree (BDT) with nodes representing tests of domain literals. The BDT is constructed incrementally over time.

Definition 2 [Relevant system description]

The system description relevant to the transition $T = \langle \sigma_1, atg, \sigma_2 \rangle$, i.e., $\mathcal{D}(T)$, is defined by signature $\Sigma(T)$ and axioms. The signature $\Sigma(T)$ is constructed to comprise the following:

1. Basic sorts of Σ that produce a non-empty intersection with $relCon(T)$.
2. All object constants of basic sorts of $\Sigma(T)$ that form the range of a static attribute.
3. The object constants of basic sorts of $\Sigma(T)$ that form the range of a fluent, or the domain of a fluent or a static, and are in $relCon(T)$.
4. Domain attributes restricted to basic sorts of $\Sigma(T)$.

Axioms of $\mathcal{D}(T)$ are those of \mathcal{D} restricted to $\Sigma(T)$. For $atg = serve(robot_1, cup_1, person_1)$ in our current example, $\mathcal{D}(T)$ does not include other robots, cups or people in the domain. It can be shown that for each transition in the transition diagram of \mathcal{D} , there is a transition in the transition diagram of $\mathcal{D}(T)$. States of $\mathcal{D}(T)$, i.e., literals comprising fluents and statics in the answer sets of the ASP program, are states in the RL formulation, and actions are ground actions of $\mathcal{D}(T)$. Also, it is possible to pre-compute or reuse some of the information used to construct $\mathcal{D}(T)$ for any given T .

Once the $\mathcal{D}(T)$ relevant to the target transition T has been identified, the RL formulation is constructed as before to compute the values of state-action combinations. The extent to which com-

Table 1. Overall control loop used by the architecture for learning and reasoning.

Input: $\Pi(\mathcal{D}, \mathcal{H})$; goal description; initial state σ_1 .
Output: Control signals for robot to execute.

```

1 planMode = true, learnType = 0
2 while true do
3     Add observations to history.
4     ComputeAnswerSets( $\Pi(\mathcal{D}, \mathcal{H})$ )
5     if planMode then
6         if existsGoal then
7             // Goal exists, consistent model, execute plan
8             if explainedObs then
9                 ExecutePlanStep()
10            else
11                // Initiate Q-RRL
12                planMode = false, learnType = 1
13            end
14        else
15            // Initiate active learning
16            planMode = false, learnType = 2
17        end
18    else
19        // Interrupt learning if needed
20        if interrupt then
21            planMode = true
22            // Continue learning
23        else if learnType == 1 then
24            ContinueRRL()
25        else if learnType == 2 then
26            if verbalCue then
27                ContinueActiveLearn()
28            else
29                ContinueActiveRRL()
30            end
31        end
32    end
33 end
    
```

puting $\mathcal{D}(T)$ reduces the search space depends on the relationships between the domain attributes and axioms. For instance, although there are several thousand static attribute combinations and more than a million object configurations in our instantiation of the RA domain, computing $\mathcal{D}(T)$

often reduces the space of attribute combinations to as few as 12 for the *serve* action. However, in other domains with complex relationships between objects, including certain variants of the classic blocks world problem in which different physical arrangements of objects with different attributes are equivalent, exploration may need to be further limited to a fraction of this restricted state space. Furthermore, Q-learning does not generalize to relationally equivalent states.

To further limit the search space and promote generalization to relationally equivalent states, we incorporate an approach inspired by the RRL-TG algorithm (Driessens & Ramon, 2003). Our approach uses experiences obtained during RL trials to construct a binary decision tree (BDT) whose nodes represent tests of domain literals—Figure 2 shows part of a BDT. Unlike the destructive branching of RRL-TG, our approach models the partial description of a state-action pair as a path to a leaf where the remaining state information is stored. When Q-value variance is reduced by adding a test at a leaf, the BDT is expanded and used to compute policies in subsequent RL episodes. To learn generic versions of the axioms, our approach explores different values of static and fluent literals. ASP-based reasoning automatically selects relevant combinations to make exploration tractable, and sampling is used if the search space is still too large. Unlike traditional RRL methods, the learned Q-values now represent values across potentially different MDPs.

After the learned Q-values converge, the system constructs axioms from the BDT, selecting the partial state-action description implied by a path to a leaf only if it is associated with a high value accrued at the leaf over multiple episodes. All subsets of the selected descriptions’ literals become candidate axioms. Because each candidate axiom could correspond to different branches of the BDT, the learner estimates each candidate’s quality by randomly drawing a number of samples without replacement from the descriptions. For each sample, it considers additional literals stored at the leaves, and revises candidates that match the sample. Each candidate axiom then records the amassed Q-value, variance, and number of training samples that influenced it.

Candidate axioms with sufficient support are *validated*, i.e., tested under simulated conditions designed to match that transition which triggered learning. Candidates that do not pass these tests are removed from further consideration. For instance, if a learned executability condition is correct, executing the action when literals in the body are true should not provide the expected outcome. Note that these tests are guaranteed to not mistakenly remove any valid axioms, although they may not eliminate all false positive candidates. The final candidates are lifted by replacing ground terms with variables, and added to the ASP program as axioms for subsequent reasoning. We refer to our RRL approach as “Q-RRL”.

Table 1 describes the overall control loop for reasoning and learning in our architecture. The baseline behavior (lines 5-17) is to plan and execute actions to achieve the given goal as long as a consistent model of history is can be computed (lines 7-9). If such a model cannot be constructed, it is attributed to an unexplained, unexpected transition, and the robot triggers Q-RRL (lines 9-12) to discover the corresponding unknown axioms (lines 20-21). If there is no active goal to be achieved, the robot triggers active learning (lines 13-16) using Q-RRL (lines 25-27) or verbal descriptions obtained from a human participant (lines 23-25) to learn previously unknown actions or axioms. When in the learning mode, the robot can be interrupted if needed (lines 18-19), e.g., to pursue a new goal. An implementation of the control loop and the components described above can be found online: <https://github.com/bmeadows/actionaxiomlearning>.

4. Experimental Setup and Results

In this section, we describe the results of experimentally evaluating four distinct hypotheses:

- **H1:** Verbal descriptions enable active learning of actions and causal laws, which serve as the foundation for further learning;
- **H2:** Q-RRL enables reliable discovery of axioms related to known or learned domain actions;
- **H3:** Learned knowledge improves the quality of plans computed for any given goal; and
- **H4:** Reliability of learning domain knowledge degrades gracefully with increase in environmental noise.

We describe execution traces in support of hypothesis *H1*, and provide quantitative results in support of hypotheses *H2*, *H3* and *H4*. Although human input is used to learn actions, it is not essential for learning the axioms; they can be learned from experiences accumulated over time.

4.1 Experimental Setup

We first: (a) experimentally set the values of parameters in Q-RRL; (b) learned model parameters used in the experimental trials; (c) defined parts of the domain knowledge to be acquired; and (d) selected performance measures and determined how their values were to be computed.

The values of the Q-RRL parameters were determined experimentally by trading off accuracy of estimated policies against processing time. For instance, learning rate and exploration preference were set to 0.1 because it resulted in a good trade-off between exploration and exploitation in the RL episodes. Also, up to 10 validation tests were conducted to evaluate candidate axioms.

We ran realistic simulation experiments using domain knowledge and statistics of action outcomes. For instance, we had the physical robots execute different actions (e.g., move between places or pick up objects) and recorded the outcomes. These statistics were used to construct T_f and R_f used in the simulation environment. Of course, both T_f and R_f were not known to the learner.

The ASP program used in the experimental trials corresponded to the RA domain of Example 1. We simulated different number of objects of different sorts, and randomly selected the location of these objects in the rooms. To simplify analysis of results, candidate axioms were constrained to have no more than two negated and two non-negated atoms of domain attributes—increases this limit also increases the computational effort of learning. We report results of experiments conducted to discover two actions (*serve* and *label*), and the corresponding axioms:

1. Serving an object to a salesperson causes it to be labelled (*causal law*);
2. A person who is not an engineer cannot be served a damaged object (*executability condition*);
3. A robot with a pneumatic arm cannot serve a brittle object (*negative affordance*);
4. A damaged object cannot be served to a person who is not an engineer, unless it is labeled (*positive affordance*);
5. An object with a brittle surface cannot be labeled by a robot (*executability condition*);
6. A damaged object cannot be labeled by a robot with a pneumatic arm (*negative affordance*);
7. Labelling a light object with a pneumatic arm causes it to be damaged (*causal law*); and

8. An object with a brittle surface cannot be labelled by a robot, unless the object is heavy and the robot has an electromagnetic arm (*positive affordance*).

The first four axioms correspond to action *serve*; the others correspond to action *label*. The action and an associated causal law were acquired from verbal descriptions. Q-RRL was then used to learn one causal law, one executability condition, one positive affordance and one negative affordance for each of the two actions. Axioms for each action can be discovered concurrently.

We used *precision* and *recall* as the key performance measures. Axioms were required to exactly match the ground truth to be counted as true positives; under-specifications (e.g., some missing literals) and over-specifications (e.g., unnecessary literals) were considered false positives. Plan quality was measured as the ability to compute a minimal plan to achieve the desired goal. Each value of these measures reported below was averaged over 1000 repetitions (e.g., for each axiom).

4.2 Execution Trace

The following execution trace supports *H1* by illustrating learning of actions and the objects those actions operate on, using verbal cues from human participants.

Execution Example 1 [*Learning from human input*]

Suppose the robot in the RA domain does not know actions *label* and *serve*, or the related axioms. For each action, the agent receives five grammatically-correct descriptions of the action being applied by another robot; these statements uphold our assumptions but otherwise vary arbitrarily. First consider the action *label*:

- The learner receives “A robot is labeling the lightweight cup” and the observation $labeled(cup_1)$. It parses the statement, matches it to the domain, lifts it to $label(\#robot, \#cup)$, and infers:

$$label(R, B) \text{ causes } labeled(B)$$

- Next, the learner receives “Robot labeled one computer”, and $labeled(comp_1)$. It learns the signature $label(\#robot, \#computer)$ and generalizes over the learned signatures to obtain $label(\#robot, \#object)$.
- Further input descriptions are automatically reconciled either when specific sorts are subsumed by more general ones, e.g., when it learns from “The pneumatic robot labels the light breakable cup”, or parsing the input results in an exact match for the action description, as in “Next the robot labeled the hard, hefty item”.

Next, in the context of the *serve* action:

- The learner observes $in_hand(p_1, book_1)$ to be true and receives “A robot serves a manual to the manager”. It produces the action description $serve(\#robot, \#book, \#person)$ and extracts the causal law:

$$serve(R, O, P) \text{ causes } in_hand(P, O)$$

- Next, the learner observes $in_hand(p_0, cup_1)$ to be true and receives “The pneumatic robot is serving the breakable cup to the clerical person over there”. Generalizing over the two examples results in $serve(\#robot, \#object, \#person)$. The remaining sentences, “Robot serves ledger

to clerical person” and “A robot served a lightweight cup to an expert”, fit the inferred structures and do not change them.

For both actions, two examples were sufficient to reach the required level of generality to model the action and an initial causal law. A key advantage of learning from verbal cues is that only a small number of examples are needed to learn the actions and the objects that they operate on. This is especially useful when actions have irreversible effects. The disadvantage is that humans are expected to provide correct descriptions of the behaviors they observe, although the robot can revise any incorrect information learned and included in the ASP program.

During the evaluation of *H1*, we considered learning incorrect knowledge, or failing to learn, from the human input as failures. The learner will fail if given incorrect inputs, e.g., “A robot serves a manual to the manager” with the observation of $in_hand(rob_1, book_1)$. The learner will learn incorrect knowledge if the POS tagger provides incorrect tags when the input strings are ambiguous, e.g., “The robot serves the cat food”. Also, a mismatch between verbal cues and WordNet synsets will result in failures, e.g., “lighter” will not map to “light” if the string “The robot is attaching a label to the lighter cup” is provided as input.

Despite these limitations, the distributed representation of knowledge supports some key capabilities. First, it simplifies inference and information reuse, e.g., if a cup has a graspable handle, this relation also holds true for other objects with handles. If an affordance prevents the robot from picking up a heavy object, this information may be used to infer that it cannot open a large window. This relates to research in psychology which indicates that humans can judge action capabilities of others without actually observing them perform the target actions (Ramenzoni et al., 2010). Second, it becomes possible to respond efficiently to queries that require consolidation of knowledge across different attributes of objects or robots, and to develop composite affordance relations, e.g., a hammer may afford an “affix objects” action in the context of a specific agent because the handle affords a pickup action and the hammer affords a swing action, for the agent. Finally, learning from verbal descriptions can be used to provide more meaningful explanations of decisions.

4.3 Quantitative Evaluation

Next we describe the results of evaluating hypotheses *H2* and *H3* in a simulated environment.

H2: Q-RRL enables reliable discovery of axioms. We explored the ability of Q-RRL to support the learning of previously unknown axioms related to a known or newly learned action. Table 2 summarizes results averaged over the four axioms for each action, showing that Q-RRL attains high recall and precision, especially after the candidate axioms are validated. The accuracy of discovering the axioms for *serve* is lower than that for *label* because the action *serve* is more complex, i.e., it has more arguments than the action *label*. There were very few differences in the values of performance measures for causal laws, executability conditions and negative affordances. The recall and precision measures were a little lower for positive affordances because axioms corresponding to positive affordances are more complex: they add context to an executability condition to make an action applicable which would be inapplicable without that context. Furthermore, note that human input is not essential for the learning of axioms with Q-RRL – a robot can learn the axioms from experiences accumulated over time through active observation or reactive action execution.

Table 2. Accuracy when Q-RRL was used to discover multiple axioms corresponding to two specific actions: *label* and *serve*. Q-RRL provides high recall and precision, especially after candidate axioms are validated.

Action	Recall	Precision	Precision (validated)
label	0.92	0.82	0.96
serve	0.88	0.70	0.95

H3: Learning improves plan quality. Next, we explored the effects of the discovered axioms on the system’s ability to generate minimal plans that provide the correct outcome. For each axiom of each target action, we conducted 1000 paired ASP-based planning trials with and without the corresponding target axiom in the system description. Each trial used a randomized scenario that required the target action to achieve the goal. We found that adding the learned executability conditions or negative affordances to the ASP program resulted in 13% (*serve*) or 23% (*label*) fewer plans. Executability conditions and negative affordances preclude certain actions in some contexts, i.e., these results indicate that the acquired knowledge improved the quality of the computed plans by eliminating plans that were incorrect or suboptimal. In a similar manner, learning and including knowledge of previously unknown positive affordances in the ASP program resulted in 17% (*serve*) or 23% (*label*) more plans. Recall that knowledge of positive affordances enables the consideration of previously unknown state transitions. Since planning (in our architecture) includes minimality constraints, these results imply that the robot can choose from a broader set of options during execution, potentially using additional heuristics or costs associated with different actions. For instance, the robot may have figured out that it can push a particular object with its arm, and may chose to use this action to move the object because pushing the object consumes less energy than lifting and carrying the object. Additional trials in which we included or removed all the learnable axioms collectively resulted in a difference of 19% (*serve*) or 58% (*label*) in the plans found. We verified that all the plans computed after including the target axioms were minimal and correct, i.e., of the best quality possible. These results indicate that learning axioms eliminates incorrect or suboptimal plans, and helps compute plans that would not be computed in the absence of the learned axioms.

In the paired trials that included or excluded the causal laws extracted from the verbal cues, there was no measurable difference in the number of plans found. This is expected; a causal law for *serve* produces outcomes which impact the applicability of other actions, and similarly for *label*. A similar situation will arise for any scenario in which the plan produced does not repeat the action influenced by the causal law. In alternative runs that involved planning for a random goal, we observed that the presence or absence of causal laws had an impact on the number of plans found.

H4: Performance degrades gracefully with environmental noise. To measure performance degradation as a function of environmental noise, we introduced a probability with which the learner could encounter simulated actuator noise with an action during learning. This noise took the form of a fixed chance for an action’s outcomes to include the removal or addition of a single random fluent of the state. We varied noise from 0 – 20%. At each discrete step in this interval, we performed 1000 trials for each target axiom. We observed a steady decline in precision and recall of axiom discovery with the increase in noise, as summarized in Table 3.

Table 3. Accuracy in the presence of actuator noise. Recall and precision decrease smoothly with the increase in the probability of encountering an actuator error while executing an action.

Noise	Recall	Precision	Precision (validated)
0%	0.89	0.76	0.95
2%	0.85	0.59	0.81
4%	0.81	0.50	0.73
6%	0.77	0.44	0.66
8%	0.72	0.37	0.59
10%	0.69	0.33	0.53
12%	0.67	0.30	0.48
14%	0.64	0.27	0.44
16%	0.61	0.24	0.40
18%	0.57	0.22	0.36
20%	0.56	0.20	0.34

Many of the false positives acquired were overly-specific variants of correct axioms. These are common when the learner observes a target axiom to fail due to noise. Without noise, 26% of false positives are over-specifications, but this increases to 41% at 10% noise, and 47% at 20% noise. Validation testing cannot adequately remove all these over-specifications, as they are in a sense correct: after validation, 78–80% of the observed errors are over-specifications. Also, consistently 1–2% of false positives are re-wordings (i.e., logical equivalences) of some target axiom. An evaluation scheme assigning a weaker penalty for over-specifications and logical equivalences in the discovered axioms would result in a slower decrease in precision with increasing levels of noise.

Our evaluation also included other findings that we did not elevate to the level of hypotheses. For instance, we found that using the ASP-based inference to guide learning makes learning significantly more efficient. We also observed that the relational representation significantly speeds up reinforcement learning in comparison with not using the relational representation.

4.4 Uncertainty, Scaling, and Computational Effort

Recall that uncertainty in perception (i.e., symbol or observation) has been abstracted away in the architecture described above. This simplification helped us focus on the interplay between representation, reasoning, and learning, but there is uncertainty in both perception and actuation when robots are used in dynamic, partially observable domains. However, since the underlying knowledge representation and reasoning architecture supports non-monotonic logical reasoning and probabilistic reasoning (Sridharan et al., 2017a), we can easily introduce probabilistic models of the uncertainty in perception. Introducing such models will change the mathematical formulation of learning from repeated trials and delayed rewards (as in Section 3.2) from an MDP to a POMDP. Although such a formulation will significantly increase the computational effort involved in learning domain knowledge, we hypothesize that learning will remain tractable because our approach uses ASP-based reasoning to automatically limit learning to the relevant parts of the domain. We leave further exploration of this idea as a direction for future research.

Next, we consider our architecture’s scalability to more complex domains, e.g., domains with many objects and complex relationships between objects. ASP-based reasoning has been shown to scale well to large domains with many objects and axioms (Erdem et al., 2016). The other key component of our architecture comprises two schemes for learning domain knowledge (see Section 3.2). The first scheme based on human verbal input involves minimal computational effort, and it can be used for real-time processing of verbal inputs. In the second scheme based on RRL, each episode of learning is not computationally expensive, but it may take multiple episodes for the values (of state action combinations) to converge and candidate axioms to be generated, especially if the knowledge being acquired is a complex relationship between different attributes. The actual learning time can thus vary significantly depending on the domain. However, our RRL approach provides an elegant means for acquiring labeled training samples associated with a measure of relative merit, in the presence of delayed action effects. Our RRL approach also incorporates a notion of relevance and exploits the underlying relational representation, e.g., ASP-based reasoning is used to guide learning. Furthermore, experimental results indicate that our RRL approach shows promise for scaling to domains with many objects and more complex relationships between objects. We leave further exploration of the scalability of our approach as a direction for future research.

5. Conclusions

This paper described an architecture for representing, reasoning with, and interactively learning actions’ names, preconditions, effects, and the objects over which these actions operate, along with associated affordances. We used Answer Set Prolog, a declarative language, to represent incomplete domain knowledge, and to perform non-monotonic logical reasoning with this knowledge for planning and diagnostics. We also demonstrated the use of ASP-based reasoning to guide the interactive learning of actions and axioms. This learning was achieved using decision-tree induction and relational reinforcement learning from observations obtained through active exploration, reactive action execution, and verbal descriptions from humans. Experimental results in the context of a simulated robot assisting humans in an indoor domain indicate that our architecture supports: (i) reliable and computationally efficient reasoning; (ii) learning of actions and axioms corresponding to different types of knowledge; (iii) improvement in the quality of plans computed for any given goal when the learned actions and axioms are included in the system description; and (iv) smooth degradation of performance with increasing levels of environmental noise.

In the future, we will explore the learning of actions and axioms in more complex domains and evaluate the architecture on physical robots, which will require the further development and integration of an existing component that reasons about perceptual uncertainty probabilistically. The long-term objective is to enable robots assisting humans to represent, reason with, and interactively revise different descriptions of incomplete domain knowledge.

Acknowledgements

This work was supported in part by the Asian Office of Aerospace Research and Development Award FA2386-16-1-4071, and Award N00014-17-1-2434 from the US Office of Naval Research. All opinions and conclusions described in this paper are those of the authors.

References

- Balduccini, M. (2007). Learning action descriptions with A-Prolog: Action language C. *Proceedings of the 2007 AAAI Spring Symposium on Logical Formalizations of Commonsense Reasoning* (pp. 13–18). Palo Alto, CA: AAAI Press.
- Balduccini, M., & Gelfond, M. (2003). Logic programs with consistency-restoring rules. *Proceedings of the 2003 AAAI Spring Symposium on Logical Formalization of Commonsense Reasoning* (pp. 9–18). Palo Alto, CA: AAAI Press.
- Bloch, M. K., & Laird, J. E. (2017). Deciding to specialize and respecialize a value function for relational reinforcement learning. *Proceedings of the Third Multi-disciplinary Conference on Reinforcement Learning and Decision Making*. Ann Arbor, MI.
- Chai, J. Y., Gao, Q., She, L., Yang, S., Saba-Sadiya, S., & Xu, G. (2018). Language to action: Towards interactive task learning with physical agents. *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence*. Stockholm, Sweden: IJCAI.
- Driessens, K., & Ramon, J. (2003). Relational instance-based regression for relational reinforcement learning. *Proceedings of the Twentieth International Conference on Machine Learning* (pp. 123–130). Washington, DC: AAAI Press.
- Erdem, E., Gelfond, M., & Leone, N. (2016). Applications of answer set programming. *AI Magazine*, 37, 53–68.
- Erdem, E., & Patoglu, V. (2012). Applications of action languages in cognitive robotics. In E. Erdem, J. Lee, Y. Lierler, & D. Pearce, (Eds.), *Correct reasoning*. Berlin: Springer-Verlag.
- Gelfond, M., & Incezan, D. (2013). Some properties of system descriptions of AL_d . *Journal of Applied Non-Classical Logics*, 23, 105–120.
- Gil, Y. (1994). Learning by experimentation: Incremental refinement of incomplete planning domains. *Proceedings of the Eleventh International Conference on Machine Learning* (pp. 87–95). New Brunswick, NJ: Morgan Kaufmann.
- Guerin, F., Kruger, N., & Kraft, D. (2013). A survey of the ontology of tool use: From sensorimotor experience to planning. *IEEE Transactions on Autonomous Mental Development*, 5, 18–45.
- Laird, J. E. (2012). *The Soar cognitive architecture*. Cambridge, MA: MIT Press.
- Laird, J. E., et al. (2017). Interactive task learning. *IEEE Intelligent Systems*, 32, 6–21.
- Langley, P., Sridharan, M., & Meadows, B. (2018). Representation, use, and acquisition of affordances in cognitive systems. *Proceedings of the 2018 AAAI Spring Symposium on Integrating Representation, Reasoning, Learning and Execution for Goal Directed Autonomy*. Stanford, CA: AAAI Press.
- Law, M., Russo, A., & Broda, K. (2018). The complexity and generality of learning answer set programs. *Artificial Intelligence*, 259, 110–146.
- Miller, G. A. (1995). WordNet: A lexical database for English. *Communications of the ACM*, 38, 39–41.

- Otero, R. P. (2003). Induction of the effects of actions by monotonic methods. *Proceedings of the Thirteenth International Conference on Inductive Logic Programming* (pp. 299–310). Szeged, Hungary: Springer.
- Ramenzoni, V. C., Davis, T. J., Riley, M. A., & Shockley, K. (2010). Perceiving action boundaries: Learning effects in perceiving maximum jumping-reach affordances. *Attention, Perception and Psychophysics*, 72, 1110–1119.
- Sarathy, V., & Scheutz, M. (2016). A logic-based computational framework for inferring cognitive affordances. *IEEE Transactions on Cognitive and Developmental Systems*, 8, 26–43.
- Scheutz, M., Schermerhorn, P., Kramer, J., & Anderson, D. (2007). First steps towards natural human-like HRI. *Autonomous Robots*, 22, 411–423.
- Simon, H. A., & Lea, G. (1974). Problem solving and rule induction: A unified view. In L. W. Gregg, (Ed.), *Knowledge and Cognition*, 15–26. Oxford, UK: Lawrence Erlbaum.
- Someya, Y. (1998). e_lemma.txt (Version 2 for WordSmith 4). From https://lexically.net/downloads/BNC_wordlists/e_lemma.txt.
- Sridharan, M., Gelfond, M., Zhang, S., & Wyatt, J. (2017a). *A refinement-based architecture for knowledge representation and reasoning in robotics*. Unpublished manuscript, <http://arxiv.org/abs/1508.03891>.
- Sridharan, M., & Meadows, B. (2017). A combined architecture for discovering affordances, causal laws, and executability conditions. *Proceedings of the Fifth International Conference on Advances in Cognitive Systems*. Troy, NY: Cognitive Systems Foundation.
- Sridharan, M., Meadows, B., & Gomez, R. (2017b). What can I not do? Towards an architecture for reasoning about and learning affordances. *Proceedings of the Twenty-Seventh International Conference on Automated Planning and Scheduling* (pp. 18–23). Pittsburgh, PA: AAAI Press.
- Sutton, R. L., & Barto, A. G. (1998). *Reinforcement learning: An introduction*. Cambridge, MA: MIT Press.
- Tadepalli, P., Givan, R., & Driessens, K. (2004). Relational reinforcement learning: An overview. *Papers from the Relational Reinforcement Learning Workshop at the Twenty-First International Conference on Machine Learning*. Banff, AB, Canada: IEEE Press.
- Toutanova, K., Klein, D., Manning, C., & Singer, Y. (2003). Feature-rich part-of-speech tagging with a cyclic dependency network. *Proceedings of the 2003 International Conference of the North American Chapter of the Association for Computational Linguistics* (pp. 252–259). Edmonton, Canada: ACM.
- Zech, P., Haller, S., Lakani, S. R., Ridge, B., Ugur, E., & Piater, J. (2017). Computational models of affordance in robotics: A taxonomy and systematic classification. *Adaptive Behavior*, 25, 235–271.
- Zhang, S., Sridharan, M., & Wyatt, J. (2015). Mixed logical inference and probabilistic planning for robots in unreliable worlds. *IEEE Transactions on Robotics*, 31, 699–713.