

GUIDE

Experiments in Socially Guided Exploration: Lessons Learned in Building Robots that Learn with and without Human Teachers

Andrea L. Thomaz^{a*†} and Cynthia Breazeal^b

^a*Georgia Institute of Technology, 85 5th Street NW, Atlanta GA 30306 USA;* ^b*MIT Media Lab, 20 Ames St., Cambridge MA 02139 USA*

(January 2008)

We present a learning system, Socially Guided Exploration, in which a social robot learns new tasks through a combination of self-exploration and social interaction. The system's motivational drives, along with social scaffolding from a human partner, bias behavior to create learning opportunities for a hierarchical Reinforcement Learning mechanism. The robot is able to learn on its own, but can flexibly take advantage of the guidance of a human teacher. We report the results of an experiment that analyzes what the robot learns on its own as compared to being taught by human subjects. We also analyze the video of these interactions to understand human teaching behavior and the social dynamics of the human-teacher/robot-learner system. With respect to learning performance, human guidance results in a task set that is significantly more focused and efficient at the tasks the human was trying to teach, while self-exploration results in a more diverse set. Analysis of human teaching behavior reveals insights of social coupling between the human teacher and robot learner, different teaching styles, strong consistency in the kinds and frequency of scaffolding acts across teachers, and nuances in the communicative intent behind positive and negative feedback.

Keywords: Human-Robot Interaction; Machine Learning; Artificial Intelligence; Social Scaffolding; Computational Models of Social Learning

1. Introduction

Enabling a human to efficiently transfer knowledge and skills to a robot has inspired decades of research. When this prior work is viewed along a guidance-exploration spectrum, an interesting dichotomy appears. Many prior systems are strongly dependent on human guidance, learning nothing without human interaction (e.g., learning by demonstration (Chernova and Veloso 2007, Atkeson and Schaal 1997, Kaiser and Dillman 1996), learning by observation (Kuniyoshi et al. 1994), learning by physical guidance (Levas and Selfridge 1984, Calinon et al. 2007), or by tutelage (Nicolescu and Matarić 2003, Lockerd-Thomaz and Breazeal 2004)). In systems such as these, the learner does little if any exploration on its own to learn tasks or skills beyond what it has observed with a human. Furthermore, the teacher often must learn how to interact with the machine and know precisely how it needs to perform the task.

Other approaches are almost entirely exploration based. For example, many prior works have given a human trainer control a reinforcement learner's reward (Blumberg et al. 2002, Kaplan et al. 2002, Saksida et al. 1998), allow a human to provide

*Corresponding author. Email: athomaz@cc.gatech.edu.

†This research was conducted at the MIT Media Lab.

advice (Clouse and Utgoff 1992, Maclin et al. 2005), or have the human tele-operate the agent during training (Smart and Kaelbling 2002). Exploration approaches have the benefit that learning does not require the human’s undivided attention. However, they often give the human trainer a very restricted role to scaffold learning, and require the human to learn how to interact with the machine.

Our research is motivated by the promise of personal robots that operate in human environments to assist people on a daily basis. Personal robots will need to be able to learn new skills and knowledge while “on the job.” Certainly, personal robots should be able to learn on their own—either discovering new skills and knowledge or mastering known skills through practice. However, personal robots must also be able to learn from members of the general public who are not familiar with the technical details of robotic systems or Machine Learning algorithms. Nevertheless, these people do bring to the table a lifetime of experience in learning from and teaching others. This is a collaborative process where the teacher guides the learner’s exploration, and the learner’s performance shapes further instruction through a large repertoire of social interactions. Therefore, personal robots should be designed to be social learners that can effectively leverage a broad repertoire of human scaffolding to be successful and efficient learners.

In sum, personal robots should be able to explore and learn on their own, but also take full advantage of a human teacher’s guidance when available.

In this paper, we present a novel learning architecture for a social robot that is inspired by theories in developmental psychology and is informed by recent advances in intrinsically motivated reinforcement learning (e.g., (Singh et al. 2005, Oudeyer and Kaplan 2004, Schmidhuber 2005)). We call this approach Socially Guided Exploration to emphasize that the robot is designed to be an intrinsically motivated learner, but its exploration and learning process can readily take advantage of a broad repertoire of a human teacher’s guidance and scaffolding.

Further, we approach this challenge from a Human-Robot Interaction perspective where we are interested in the human-teacher/robot-learner system. Hence, our analysis examines and compares the content of what is learned by the robot both when learning in isolation and with a human teacher. To do this, we conducted a human subjects experiment where we had 11 people (all previously unfamiliar with the robot) teach it to perform a number of tasks involving a puzzle box. The puzzle box is pre-programmed with a suite of behaviors such as changing the color of its lights, opening or closing its lid, or playing a song when the correct sequence of button presses, switch flips, or slider toggles is performed. We analyze our human subjects’ teaching behavior in relation to the robot’s learning behavior to understand the dynamics of this coupled social process. Our findings reveal social constraints on how people teach socially interactive robots.

The primary contribution of this article is the analysis of the social coupling between a human teacher and a robot learner. Our findings have important implications for how to design social robots that learn from everyday people.

2. Robot Platform

Our research platform is Leonardo (“Leo”), a 65 degree-of-freedom anthropomorphic robot specifically designed for human social interaction (see Figure 1). Leo has speech and vision sensory inputs and uses gestures and facial expressions for social communication (the robot does not speak yet). Leo can visually detect objects in the workspace, humans and their head pose, and hands pointing to objects. For highly accurate tracking of objects and people, we use a 10 camera VICON optical motion capture system (<http://www.vicon.com/>). The speech understanding



Figure 1. Our social robot, Leonardo, interacts with a human teacher to learn about a puzzle box.

system is based on Sphinx (Lamere et al. 2003), and has a limited grammar to facilitate accuracy.

3. Socially Guided Exploration System

In most Machine Learning systems, learning is an explicit activity. Namely, the system is designed to learn a particular thing at a particular time. In human learning, on the other hand, learning is a part of all activity. There is a motivation for learning, a drive to know more about the environment, and an ability to seek out the expertise of others. Children explore and learn on their own, but in the presence of a teacher they can take advantage of the social cues and communicative acts provided to accomplish more than they would be able to do on their own (also known as social scaffolding (L. S. Vygotsky 1978)).

A teacher often guides a learner by providing timely feedback, leading them to perform desired behaviors, and controlling the environment so the appropriate cues are salient, thereby making the learning process more effective. This is the primary inspiration for the Socially Guided Exploration system.

This section highlights the key implementation details. We first cover the representation used for objects, goals, and tasks (Sec. 3.1). Then we describe the Motivation System (Sec. 3.2) that arbitrates three high level Learning Behaviors (Sec. 3.3) that result in Task Learning and Generalization (Sec 3.4 and 3.5). Finally, we describe the behaviors the robot uses to facilitate Transparency (Sec. 3.6), and the Scaffolding mechanisms available for a human teacher to influence the process (Sec. 3.7).

3.1. *Fundamentals of Task Representation*

3.1.1. *Object Beliefs*

The Socially Guided Exploration system extends the C5M architecture (Blumberg et al. 2002)—adding capabilities for representing and learning goal-oriented tasks, self-motivated exploratory behavior, and expression/gesture capabilities to support a collaborative dialog with a human teacher. The Perception and Belief Systems of C5M are most relevant to the learning abilities described in this paper. Every time step, the robot has observations from its various sensory processes, $O = \{o_1, \dots, o_k\}$. The Perception System is a set of *percepts* $P = \{p_1, \dots, p_n\}$. Each $p \in P$ is a

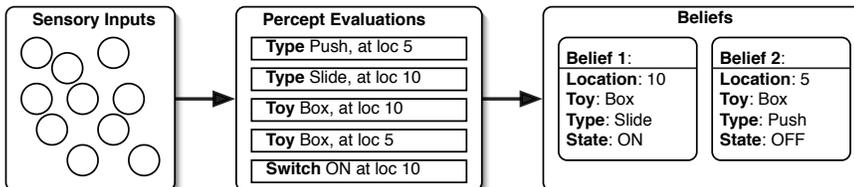


Figure 2. Sensory input is classified by percepts and then merged into discrete object representations. In this timestep, 5 percepts yield 2 object beliefs.

classification function, such that $p(o) = m$ where $m \in [0, 1]$ is a match value. The Belief System maintains the *belief* set B by integrating these percepts into discrete object representations (based on spatial relationships and various similarity metrics). Figure 2 shows a simplified example in which sensory data leads to five percepts with $m > 0$, that result in two beliefs in B . In this paper, a “state” s refers to a snapshot of the belief set B at a particular time, and S refers to the theoretical set of all possible states. Let $A = \{a_1, \dots, a_i\}$ be the set of Leo’s basic actions. For more details of the Perception and Belief Systems see (Breazeal et al. 2005).

3.1.2. Goals

Socially Guided Exploration aims to have a system learn the goal or concept of an object-oriented activity, where a goal is a particular state change. Goal beliefs are derived from the set of beliefs about objects in the world. Rather than containing a single set of percept values, a goal belief represents a desired change to an object during an action or task by grouping a belief’s percepts into *expectation* percepts (indicating an expected object feature value), and *criteria* percepts (indicating which beliefs are relevant to apply this expectation to).

The procedure for making a goal, G , given the two states, s_i and s_j is the following.¹ Create a goal belief, x , for each belief in s_i that changed over $s_i \rightarrow s_j$: $\forall b_i \in s_i$ find the corresponding² belief, $b_j \in s_j$. If there are any percepts differences between b_i and b_j then make a goal belief x in the following way: $\forall p \in b_i$ if b_j has the same value for p then add p to x as a criteria percept (i.e. add p to $crit \in x$), otherwise add the b_j value of p to x as an expectation percept (i.e. add p to $expt \in x$). Then add x to the the set of goal beliefs, G . At the end of this process, G contains a goal belief for each object that incurred any change over $s_i \rightarrow s_j$, for an example of goal inference, see Figure 3.

A goal, G , can easily be evaluated to see if it is true or complete in the current state in the following way: $\forall x \in G$, if any belief $b \in B$ matches all of the $crit \in x$, then b must also match all of the $expt \in x$.

3.1.3. Task Representation

The Socially Guided Exploration system learns and maintains a set of *Tasks*. Each $T \in Tasks$ is represented as a *Task Option Policy*.³ This name is meant to reflect its similarity to the Options approach in Reinforcement Learning (Sutton et al. 1999). Options are temporally extended actions and are represented with three constructs (I, π, β) :

$$\pi: S \times A \rightarrow [0, 1]; \text{ A policy estimating a value for } (s, a) \text{ pairs.}$$

¹This goal construct is also used in prior work, Lockerd-Thomaz and Breazeal (2004), Breazeal et al. (2005).

²“Corresponding” here refers to the fact that b_i and b_j are actually snapshots from the same belief objects in the Belief System. Recall that beliefs are collections of percept histories, thus b_i and b_j are different timeslices of the same collections of percept histories.

³In this article, ‘Task Option Policy’ is often shortened to ‘Task’ for simplicity.

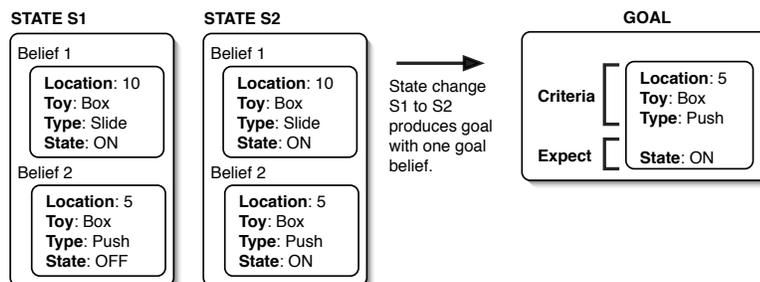


Figure 3. A simple example of creating a goal from a state change. As explained in Section 3.4, this goal inference happens each time the Novelty behavior activates and attempts a new task creation.

$\beta: S^+ \rightarrow [0, 1]$, where $S^+ \subset S$; is the termination set, all states in which the option terminates.

$I: \subseteq S$ is the initiation set, all the states in which the option can initiate.

An option can be taken from any state in I , then actions are selected according to π until the option terminates stochastically according to β .

Each Task Option Policy, $T \in Tasks$, is defined by very similar constructs (I', π', β') . To define these we use two subsets of states related to the task. Let $S_{task} \subset S$ be the states in which the task is relevant but not achieved, and $S_{goal} \subset S$ be the states in which the goal is achieved. Then, a Task Option Policy is defined by:

$\pi': S_{task} \times A \rightarrow [0, 1]$; estimates a value for (s, a) pairs in relation to achieving the task goal, G .

$\beta': S_{goal}$; represents all of the states in which this task terminates because G is true.

$I': S_{task}$; represents the initiation set—the task can be initiated in any state for which it has a policy of action.

A task can be executed (i.e., it is *relevant*) when the current state is in S_{task} . During execution, actions are chosen according to π' until the current state is in S_{goal} (with some probability of terminating early). Again, a state s achieves the goal if: $\forall x \in G$, if any belief b in s matches all the *criteria* $\in x$, then b also matches all the *expectation* $\in x$.

The following sections describe how the system learns new *Tasks* and refines their representation with experience.

3.2. Motivational Drives for Learning

Living systems work to keep certain critical features within a bounded range through a process of behavioral homeostasis (e.g., food, water, temperature). If a parameter falls out of range, the animal becomes motivated to behave in a way that brings it back into the desired range.

Recently, this concept has inspired work on internal motivations for a Reinforcement Learning (RL) agent (Singh et al. 2005, Oudeyer and Kaplan 2004, Schmidhuber 2005). These works use a measure of novelty or certainty as intrinsic reward for a controller. Thus, an action that leads to a prediction error results in rewards that encourage focus on that portion of the space. For example in (Singh et al. 2005), their ‘intrinsic’ reward is related to expected errors as represented in the transition model for a task. In (Oudeyer and Kaplan 2004), they use the first derivative of this, taking ‘progress’ to be the reduction in prediction errors of the transition model for a task.

Our approach is in a similar vein, but rather than contribute to the reward

directly, Leo’s internal motivations trigger learning behaviors that help the system arbitrate between learning a new task, practicing a learned task, and exploring the environment.

Leo’s Motivation System (based on prior work (Breazeal 2002)) is designed to guide a learning mechanism. Inspired by natural systems, it has two motivational drives, **Novelty** and **Mastery**. These are defined in detail below. Each drive has a range $[0, 1]$, initial value of 0.5, a tendency to drift to 0.0, and a drift magnitude of 0.001 (max change in a time step). The Motivation System maintains the drive values based on the status of the internal and external environment.

3.2.1. The Novelty Drive.

The **Novelty** drive is meant to maintain a general notion of the unfamiliarity of recent events. Every state transition will cause the **Novelty** drive to rise for an amount of time related to the degree of the change, d_{chg} , based on the event’s frequency. In particular, the degree of change between state s_1 and state s_2 is inversely related to how often this state change has been seen:

$$d_{chg}(s_1, s_2) = \frac{1}{frequency(s_1, s_2)}$$

An event causes the **Novelty** drive to drift towards its maximum value for a period of time, $t = d_{chg}(s_1, s_2)t_{max}$. The maximum effect time, t_{max} , is 30 seconds.¹

3.2.2. The Mastery Drive.

The **Mastery** Drive reflects the current system confidence of the *Tasks* set. **Mastery** is the average confidence of the tasks that are relevant in the current state, s (i.e., tasks that can be initiated from s). A task’s confidence, C , is the number of successful attempts, reaching a state of the world in which the task’s goal evaluates to true (see Sec. 3.1), over the total task attempts made: $C = \frac{successes}{total}$. Thus, the following formula is used to calculate **Mastery** for a particular time step: Let X be the subset of all $T \in Tasks$ for which the current state, s , is in the initiation set S_{task} of T . **Mastery** is the average confidence of all the tasks in X .

3.3. Learning Behaviors for Motivational & Social Contexts

The Task Learning Action Group is the piece of the Socially Guided Exploration system responsible for identifying and responding to learning opportunities in the environment. It maintains the set of known *Tasks*, and has three competing learning behaviors that respond to social and motivational learning contexts. Figure 4 is an overview of the behaviors and their internal/external triggering contexts.

3.3.1. The Novelty Behavior.

One purpose of the novelty drive is to encourage the system to better understand new events, expanding the *Tasks* set. Thus, a significant rise in the novelty drive makes the Novelty behavior available for activation. Additionally, this behavior may be activated due to a social context, when the human points out an event (e.g., “Look Leo, it’s **TaskName-X**.”). Once activated, the Novelty behavior tries to create a new task. It makes a goal representation of the most recent state transition (s_1, a, s_2), and if there is not a $T \in Tasks$ with this goal, then a new task is created. Task creation, expansion, and generalization are covered in Sec. 3.4 and 3.5.

¹This value was set empirically, as a reasonable amount of time to provide activation to the learning behaviors described in the next section.

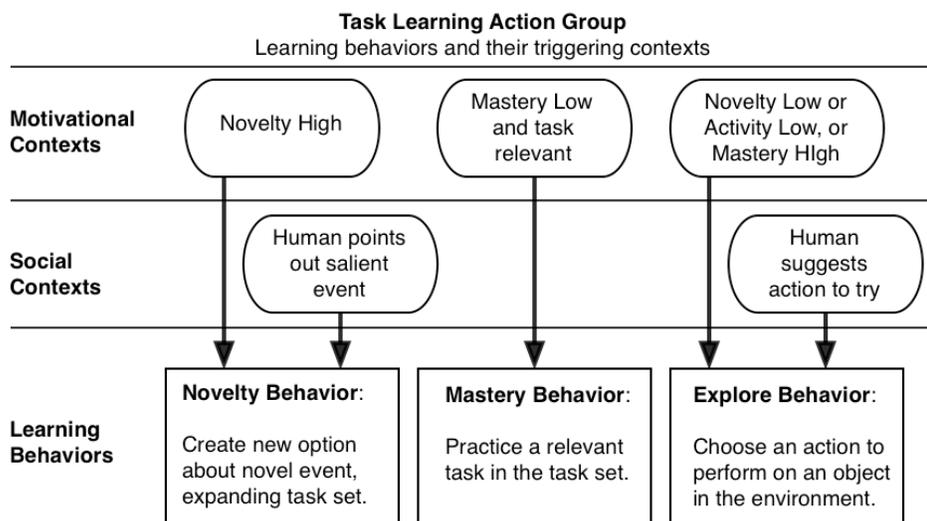


Figure 4. The three learning behaviors and their social/motivational contexts.

3.3.2. The Mastery Behavior.

The purpose of the Mastery Drive is to cause the system to become confident in the environment, fleshing out the representations in the *Tasks* set. When the Mastery Drive is low and any tasks are relevant in the current state, the **Mastery Behavior** may be activated. This behavior randomly selects a relevant task, executes it, and updates the confidence based on success in reaching the goal.

3.3.3. The Explore Behavior.

Both motivational drives also work to encourage exploration. The **Explore Behavior** becomes available when **Novelty** is low, encouraging the system to seek out the unexpected. Exploration is also triggered when **Mastery** is high. Even if a known task is relevant, the system is biased to try to expand the *Tasks* set once confidence is high. Additionally, social interaction can trigger the **Explore Behavior** — for example, if the human suggests an action (e.g., “Leo, try to Act-X the Obj-Y.”). When the **Explore Behavior** is activated, it first tries to do any human-suggested action if possible. Otherwise, the **Explore Behavior** selects from the actions it can do in the current state.¹ Once the action is completed, if it was a human-suggested action, the robot’s attention is biased to look to the human in order to acknowledge the suggested action and provide the human with an opportunity for feedback.

3.4. Learning new Tasks

The Socially Guided Exploration system learns a new Task Option Policy by creating a goal G about a state change and refining S_{task} , G , and π' over time through experience.

When the **Novelty Behavior** is activated, a potential goal state, G , is made from the most recent state change, (s_1, a, s_2) , as described in Sec. 3.1. If there does not exist a $T \in Tasks$ with goal G , then a new Task Option Policy, T_{new} , is created. The S_{task} of T_{new} is initialized with the state s_1 , and π' is initialized with

¹The exploration action selection mechanism is very simple. It makes a random selection among available actions that have been executed the least. Thus, the system will choose random actions, but will try all actions the same number of times, all things being equal.

default values $q = .1$ for all actions from s_1 . Then, the system takes into account the experience of (s_1, a, s_2) , and (s_1, a) gets a higher value since G is true in s_2 .

Each $T \in Tasks$ can learn and expand from every experience (also referred to as intra-option learning, Sutton et al. (1998)). Every action is an experience, (s_1, a, s_2) ; and each $T \in Tasks$ has the opportunity to extend its set S_{task} and update its π' based on this experience. To update π' , the system estimates the reward function based on the task's goal: $r = 1$ if G is true in s_2 , otherwise $r = 0$.

- Extend: For each $T \in Tasks$, if a task's goal is true in s_2 or the task's initiation set S_{task} contains s_2 , then s_1 should be added to the S_{task} .
- Update: For each $T \in Tasks$, if s_1 is in the initiation set S_{task} then update the value of (s_1, a) in the π' : $Q[s_1, a] = Q[s_1, a] + \alpha(r + \gamma \max_a(Q[s_2, a]) - Q[s_1, a])$.

3.5. Task Generalization

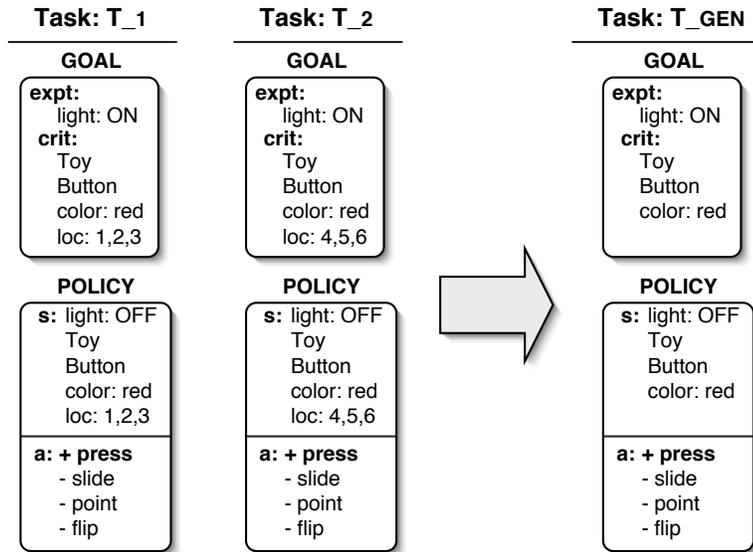
In addition to expanding initiation sets and updating value estimates for tasks, the system tries to generalize tasks over time. While the Novelty behavior creates new specific task representations about novel state changes, the system needs a generalization mechanism in place to avoid an explosion of very specific tasks being created. Thus, every time a new task is added to the $Tasks$ set, the learning mechanism works to generalize both the state representations in S_{task} and the goal representation G for all $T \in Tasks$.

Given two different tasks T_1 and T_2 , the generalization mechanism attempts to combine them into a more general task T_{gen} . For example, if T_1 has the goal of turning ON a red button in location (1,2,3), and T_2 has the goal of turning ON a red button in location (4,5,6), then T_{gen} would have the goal of turning ON a red button without a location feature. When a feature is generalized from the goal, the system also tries to generalize the states in S_{task} , letting the task ignore that feature. Thus, T_{gen} can initiate with a red button in any location and any state with a red button ON achieves its goal.

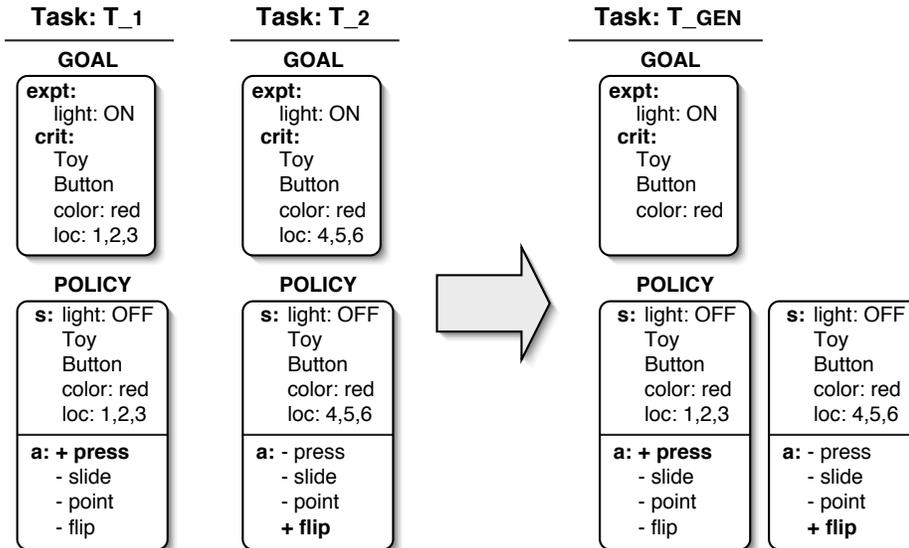
This generalization is attempted each time a T_{new} is added to $Tasks$. If there exist two tasks T_1 and T_2 with similar goal states, then the system makes a general version of this task. Two goals are similar if they differ by no more than four percepts.¹ In generalizing S_{task} and G for all $T \in Tasks$, the generalization mechanism expands the portion of the state space in which tasks can be initiated or considered achieved. This results in an efficient representation, as the system continually makes the state space representations more compact. Additionally, it is a goal-oriented approach to domain transfer, as the system is continually refining the *context* and the *goal* aspects of the activity representation.

In our red button example, the two tasks are similar since their expectations are the same, $expt = \{ON\}$, and their criteria differ only by the location feature. A new task is made with a goal that does not include location: $G_{gen} = \{expt = \{ON\}; crit = \{object, red, button, \dots\}\}$. If the policies of the two tasks are similar, for example to do the press action in the state $\{b_1 = \{object, red, button, loc = (x, y, z), \dots\}\}$, then the new task will generalize location from all of S_{task} (see Figure 5(a)). On the other hand, if T_1 has the policy of doing the press action in state $\{b_1 = \{object, red, button, loc = (1, 2, 3), \dots\}\}$, and T_2 has the policy of doing the flip action in state $\{b_1 = \{object, red, button, loc = (4, 5, 6), \dots\}\}$, then the generalized task policy will maintain that in $loc(1,2,3)$ a red button should be pressed to make

¹Four is somewhat arbitrary, chosen empirically as a good balance between under and over utilization of the generalization mechanism. This number is obviously sensitive to the task domain and feature space. Determining this parameter automatically is left as future work.



(a) Task T_1 and T_2 have similar goals, to turn the red button ON. So a general task T_{gen} is made with the generalized G , S_{task} , and π' , that no longer include the location feature.



(b) Task T_1 and T_2 have similar goals, to turn the red button ON. So a general task T_{gen} is made with the generalized G . But they have different ways of achieving this goal, so the S_{task} and π' are not generalized, but include the S_{task} and π' from both T_1 and T_2 .

Figure 5. Between-policy generalization, a simplified example: Fig. 5(a) shows the generalization for the example where the two tasks have similar goals and action policies. Fig. 5(b) shows the example where they have similar goals but different action policies.

it ON and in loc(4,5,6) a red button should be flipped (see Figure 5(b)).

3.6. Transparency Mechanisms

Leo has several expressive skills contributing to the robot’s effectiveness as a social learner. Many are designed around theories of human joint activity (Clark 1996). For example, consider principles of grounding. In general, humans look for evidence that their action has succeeded. This extends to joint activity where the ability to establish a mutual belief that a joint activity has succeeded is fundamental to a successful collaborative activity.

Table 1. Social Cues for Transparency in a Socially Guided Exploration

Context	Leo's Behavior	Intention
Human points to object	Looks at Object	Shows Object of Attention
Human present in workspace	Gaze follows human	Shows social engagement
Executing an Action	Looks at Object	Shows Object of Attention
Human says: "Look Leo, it's Task-X"	Subtle Head Nod and Happy facial expression	Confirms goal state of task-X
Human says: "Try to Act-Y the Obj-Z"	Look to human if suggestion is taken	Acknowledge partner's suggestion
Speech didn't parse, Unknown object request, Label without point	Confusion gesture	Communicates problem
Unconfident task execution	Glances to human more	Conveys uncertainty
Task is done, and human says: "Good!"	Nods head	Positive feedback for current option
Human asks a yes/no question	Nod/Shake	Communicates knowledge/ability
Intermittent Body motion	Eye blinks, Gaze shifts,	Conveys awareness and aliveness
Novel event	Surprise (raise brows/ears, open mouth)	Task being created.
Mastery triggers a task execution	Concentration (brows/ears down)	A known task is being tried
Completion of successful task attempt	Happy (open mouth, raised ears)	Expectation met
Completion of failed task attempt	Sad (closed mouth, ears down)	Expectation broken
Feedback from Human	Happy/Sad	Acknowledges feedback



Figure 6. Leo can use several facial poses to express internal learning state.

Table 1 highlights many of the social cues that Leo uses to facilitate the collaborative activity of learning. Eye gaze establishes joint attention, reassuring the teacher that the robot is attending to the right object. Subtle nods acknowledge task stages, e.g., confirming when the teacher labels a task goal.

Additionally, Leo uses its face for subtle expressions about the learning state. The robot's facial expression shifts to a particular pose for fleeting moments (2-3 seconds), indicating a state that pertains to its internal learning process, and then returns to a neutral pose. The expressions are chosen to communicate information to the human partner. They are inspired by research showing that different facial action units communicate specific meanings (Smith and Scott 1997) (Figure 6). For example, raised eyebrows and wide eyes indicate heightened attention, which is the desired communicative intent with Leo's surprised expression. This approach results in a dynamic and informative facial behavior.

Leonardo also communicates various learning contexts to the human partner with its facial expressions (Table 1). When the **Novelty Behavior** is triggered, a fleeting surprised expression lets the human know that a task is being created. When the **Mastery Behavior** causes a task to be practiced, Leo makes a concentrated facial expression and later a happy/sad expression upon the success/failure of the

attempt. Throughout, if the human gives positive or negative feedback, Leo makes a happy or sad expression to acknowledge this feedback. When the human labels a goal state, Leonardo makes a happy expression and a head nod to acknowledge the labeling.

3.7. *Scaffolding Mechanisms*

The goal of our approach is for a robot learner to strike a balance between learning on its own and benefiting from the social environment. The following are social scaffolding mechanisms at work on the Leonardo platform to enable Socially Guided Exploration.

Social attention: The attention of the robot is directed in ways that are intuitive for the human. Attention responds to socially salient stimuli and stimuli that are relevant to the current task. The robot tracks the pointing gestures and head pose of a human partner, which contribute to the saliency of objects and their likelihood for attention direction. For details on the robot’s social attention system see (Thomaz et al. 2005).

Guidance: Throughout the interaction, the human can suggest actions for Leo to try. The human’s request is treated as a suggestion rather than an interrupt. The suggestion increases the likelihood that the **Explore Behavior** will trigger, but there is still some probability that Leo will decide to practice a relevant task or learn about a novel state change.

Recognizing goal states: Leo creates task representations of novelties in the environment. The human can facilitate this process by pointing out goal states with a variety of speech utterances (e.g., “Look Leo, it’s X”). This serves to increase the likelihood that the **Novelty Behavior** will trigger, creating a task with the label “X.”

Testing knowledge: Once the human has labeled a goal state “X,” then they can help the robot decide when it is appropriate to practice the new task by suggesting a task attempt (“Leo, try to “X”). This serves to increase the likelihood that the **Mastery Behavior** will trigger, and attempt the task with the label “X.”

Environmental structure: An implicit contribution of the human teacher is their ability to physically structure the learning environment, highlighting salient elements. They draw the robot learning system into new generalizations, link old information to new situations, and point out when a learned task is relevant in the current situation.

4. Experiment

To evaluate our Socially Guided Exploration system, we conducted an experiment where human subjects interacted with the Leonardo robot. This experiment addresses two questions about the system.

- (1) What differences are seen, on average, when the robot learns by itself versus when it learns with the guidance of a human teacher?
- (2) How do human teachers make use of the scaffolding mechanisms available in the Socially Guided Exploration system?

To study these questions, we solicited volunteers from the campus community, and had 11 participants complete the experiment over the course of two days (5 male, 6 female). None of the participants had interacted with the Leonardo robot previously. Due to corrupted log files for two subjects, we are only able to use data

from 9 of the subjects in our first analysis that depends on those log files. However, for the subsequent video analysis, we use the data from all 11 subjects. In this section we describe the experiment, and in Section 5 we present two sets of results from the experiment addressing the two questions raised above.

4.1. *Experimental Scenario*

The experimental scenario is a shared workspace where Leo has a puzzle box (see Figure 1). The puzzle box has three inputs (a switch, a slider, and a button), a lid that can open and close by activating an internal motor, five colored LEDs, and sound output. The box can be programmed with specific behaviors in response to actions on the input devices (e.g., the actions required to open the lid, or turn a colored LED on, etc.).

Leo has five primitive manual actions it can apply to the box (Button-Press, Slider-Left, Slider-Right, Switch-Left, Switch-Right), but no initial knowledge about the effects of these actions on the puzzle box. Leo uses the Socially Guided Exploration mechanism to build a *Tasks* set about the puzzle box. In our experiment, the puzzle box is pre-programmed with the following input-output behavior:

- Pressing the button toggles through the five LED colors: white, red, yellow, green, and blue.
- If both the slider and the switch are flipped to the left when the color is white, then the box lid opens.
- If the slider and switch are flipped to the right when the color is yellow, then the box lid closes.
- If the lid is open and the color changes to blue, then the box will play a song.

4.2. *Instructions to Human Subjects*

Subjects are shown the functionality of the puzzle box and told that their goal is to help Leo learn about it. They are told the robot is able to do some simple actions on the toy puzzle box, and that once turned on, Leo will start exploring what it can do with the box. Then the scaffolding mechanisms are explained. They are told they can help Leo learn tasks by making action suggestions, by naming aspects of the box, and by testing that these named aspects have been learned. The subjects were told that Leo understands the following kinds of utterances:

- “Leo, try to...[press the button, move the slider/switch left/right, move the switch left/right].”
- “Look Leo, It’s...[Open, Closed, A Song, Blue, White, Green, Red, Yellow].”
- “Leo, Try to make it...[Open, Closed, Play a Song, Blue, White, Green, Red, Yellow].”
- “Good Leo”, “Good job”, “Well done”, “No”, “Not quite.”

Finally, it is explained that their goal in this interaction is to make sure that Leo learns to do three things in particular:

- T_{Blue} —Make the light blue;
- T_{Open} —Make the lid open;
- T_{Song} —Make the song play.

After receiving the instructions, subjects were given a chance to manipulate the puzzle box, before Leo was turned on, until they felt they understood the functionality. Then the learning session would begin. Each subject interacted with Leo for 20-30 minutes.

4.3. Metrics Collected

During both the self learning and guided learning sessions, we logged several measures to characterize the learning process: Actions performed by the robot, speech utterances from the human partner, each time the box transitions into one of the goal states, each time a new task is added to the *Tasks* set or modified through generalization. Finally, the resulting *Tasks* set is saved at the end of the learning session.

After each learning session, additional measures are collected in simulation about the efficacy of the learned *Tasks* sets. In this analysis, we evaluated each *Tasks* set, measuring their ability to achieve each of the experimental goals from a test suite of five initial states. Each experiment goal has a different test suite of five initial states: three of which are very close to the goal (1 or 2 actions required), two of which are farther away (more than 2 actions required to achieve the goal). For each of the learned *Tasks* sets, we measure whether or not an agent with this set of task representations is able to reach the goal from the initial conditions, and if so the number of actions needed to reach the goal from each of the test states.

We also wanted to measure how ‘related’ each of the *Tasks* sets were to the experimental goals. Thus, for each *Tasks* set we record the number of tasks related to each experiment goal. For example, a task is considered related to T_{Blue} if the blue light ‘ON’ is in the expectation portion of the task’s goal state.

5. Experimental Results

5.1. Analysis of Guided Exploration versus Self Exploration

This first set of analyses examines how the teacher’s social scaffolding influenced the robot’s learning process. We compare data from the learning sessions in two conditions:

- GUIDED: The robot learns with a human teacher. As mentioned above, we have data from 9 participants in this condition.
- SELF: The robot learns by itself. For this condition, we collected data from 10 sessions of the Leonardo robot learning alone in the same environment.

All 9 participants succeeded in getting the robot to reach the T_{Blue} and T_{Open} tasks, but none of the participants focused on the more complex T_{Song} . Thus our analysis is focused on the T_{Blue} and T_{Open} tasks. Everyone taught the T_{Blue} task first, and there was an average of 9 actions between first encountering the T_{Blue} and T_{Open} goals.

The differences between the Self Exploration and Socially Guided Exploration cases are summarized in Table 2. We found that the human teacher is able to guide the robot to the desired goal states faster than it discovers them on its own. This is seen in the difference between groups in the number of actions to the first encounter of any of the experiment goal states. The average for GUIDE, 3.56, is significantly less than the average for the SELF condition, 11.2. Thus, people were able to use the social scaffolding mechanisms to focus the robot on aspects of the environment that they wanted it to learn.

This is also supported by qualities of the resulting *Tasks* set that is learned. In the GUIDE condition, the resulting *Tasks* sets were more related to the experiment goals (i.e., T_{Blue} , T_{Open} or T_{Song} is true in a task’s goal state). We see a significant difference in both the number of tasks related to T_{Blue} and T_{Open} .

Also, we found that the Socially Guided Exploration case learns a better task set for achieving the experiment goals. The post-evaluation of the learned task sets

Table 2. Summary of differences found between Self Exploration and Socially Guided Exploration.

Mesure	Mean SELF	Mean GUIDE	T-test Results
# actions to reach first goal in learning session	11.2	3.56	$t(19) = 2.11$; $p < .05$
Size of resulting <i>Tasks</i> set	10.4	7.55	$t(19) = 7.18$; $p < .001$
# tasks for T_{Blue}	0.833	1.333	$t(19) = -2.58$; $p < .01$
# tasks for T_{Open}	1	1.77	$t(19) = -1.83$; $p < .05$
# Init States can reach T_{Open} in post-experiment	.58	1.56	$t(19) = -2.88$; $p < .01$
# actions to reach T_{Blue} in post-experiment	2.66	1.69	$t(19) = 2.19$; $p < .05$

found differences in the generality of the learned tasks. The average number of test states that the GUIDE condition sets could reach the T_{Open} goal, 1.56, was significantly better than the average in the SELF condition, 0.58. And though we didn't find this particular difference for the T_{Blue} goal, we do see that the GUIDE condition is significantly faster at achieving T_{Blue} in the post analysis than the SELF condition, 1.69 versus 2.66.

Thus, we can conclude that the nature of what is learned between self-learning and guided learning sessions is significantly different. The human teachers were able to utilize the guidance mechanisms available to help the robot learn task sets that are better at achieving the designated experimental goals.

5.2. Analysis of Human Scaffolding Behavior

5.2.1. Video Data

Having learned about how human scaffolding changes the nature of what is learned during an exploration session, our next set of analyses focuses on understanding how people used the scaffolding mechanisms provided. We have video from each of the learning sessions, and we coded a transcript from each video that summarizes the following:

For each of the human's utterances, we coded the *type* of scaffolding as one of the following:

- Action suggestion; (e.g., "Try to press the button.")
- Task label; (e.g., "Look Leo, it's Blue.")
- Task test; (e.g., "Leo, try to make it Blue.")
- Positive or negative feedback.

For each of the human's utterances, we coded its *context* as one or more of the following:

- Context 1: The person waited for the robot to make eye contact before they made the utterance.
- Context 2: The utterance happened after the robot made a facial expression.
- Context 3: The utterance happened during the robot's action.

The transcript also logs each action and emotional expression made by the robot. Though the camera was placed such that it was difficult to see every small facial expression made by the human, we logged all visible and audible emotional expressions (smiles, laughs, etc.).

Table 3. Relative frequency of scaffolding use.

Scaffolding Type	Average	Variance
Action suggestions:	0.85	0.038
Task labels:	0.36	0.022
Task tests:	0.17	0.007
Positive feedback:	0.31	0.027
Negative feedback:	0.11	0.001

Table 4. Context of scaffolding utterances.

Scaffolding Type	Context 1: eye contact		Context 2: expression		Context 3: during action	
	Average	Variance	Average	Variance	Average	Variance
Action suggestions:	0.892	0.011	0.084	0.003	0.075	0.007
Task labels:	0.879	0.012	0.051	0.003	0.136	0.027
Task tests:	0.971	0.006	0.049	0.008	0.000	0.000
Positive feedback:	0.469	0.049	0.050	0.005	0.582	0.059
Negative feedback:	0.538	0.155	0.000	0.000	0.631	0.091

5.2.2. Results

Table 3 summarizes the frequency with which each of the subjects used each type of utterance. After normalizing by length of the learning session, we calculate the average frequency (i.e., number of utterances per action) for each scaffolding type. Interestingly, we found there was little variance in these frequencies across the participants.

We also looked at whether each type of scaffolding utterance happened in a particular context. Table 4 shows the average percentages for each of the three contexts for each type of scaffolding utterance. We see that nearly all (97%) of the task tests happen after eye contact (Context 1). Action suggestions are similar — 89% are in Context 1.

With task labels, 88% happen in Context 1. However, there is some tendency (14%) toward Context 3 where people label a goal state during Leo's action, before Leo looks up. Feedback is the most varied in terms of context, and like goal labeling it occurs in both Contexts 1 and 3, as seen in Table 4. But on an individual level people are more diverse in their behavior. Two of the 11 people issued most of their feedback (more than 65%) in Context 1; six people did most of their feedback in Context 3; and three people split their utterances nearly 50/50 between Contexts 1 and 3. We were interested to find this division of context; our hypothesis for future work is that a goal label or feedback utterance given in Context 1 carries a different meaning than one given in Context 3.

In addition to the context of feedback, we looked at the relative amounts of positive and negative utterances (Figure 7). Again we have diversity among the 11 participants. We see that 3 people gave only positive feedback, 6 people had a positive bias to their feedback, and 2 people had fairly even amounts of positive and negative feedback. It is interesting that none of our subjects had a negative feedback bias.

The final question we looked at in the video analysis was “How often do people mirror the expressions/affect that Leo displayed?” Table 5 summarizes this data. There was a fairly wide range of matching behavior. Two people never did any visible or audible mirroring, while one person matched nearly 50% of Leo's expressions. On average people matched Leo's emotional expression 24% of the time.

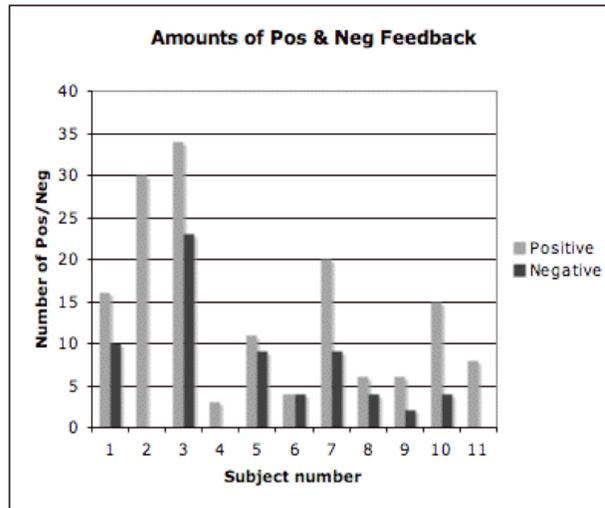


Figure 7. Relative amounts of positive and negative feedback issued by our human subjects.

Table 5. Summary of each person’s mirroring behavior.

Subject	Total number of Expressions by Leo	Number of expressions that were matched	ratio
1	10	0	0.000
2	16	2	0.125
3	13	3	0.231
4	4	0	0.000
5	14	5	0.357
6	17	8	0.471
7	12	4	0.333
8	7	1	0.143
9	10	4	0.400
10	6	2	0.333
11	9	2	0.222
Average			0.238

6. Discussion

In designing robotic agents that learn new skills and tasks “on the job” from everyday people, we recognize that the average person is not familiar with machine learning techniques, but they are intimately familiar with various forms of social learning (e.g., tutelage, imitation, etc.). This raises important research questions. For instance, “How do we design robots that learn effectively from human guidance?”; “What does human teaching behavior look like when interacting with a social robot?”; “How do we design robot learners to support human teaching behavior?” Also, there still remains the more traditional robot learning question, “How can robots be effective learners when a human teacher is not available?”

In this work, we recognize that a social learner needs both self exploration and guidance, and we bring these together in one learning system. Motivations drive exploration of the environment and arbitrate the creation of goal-oriented tasks about novel events. A human partner can influence learning through typical scaffolding acts such as directing attention, suggesting actions, highlighting and labeling goal states as interesting states to learn to achieve, testing task knowledge, and issuing positive/negative feedback.

Our experiments show that the Socially Guided Exploration mechanism is successful in allowing non-expert human teachers to guide the robot’s learning process. People were able to focus the robot’s learning to particular goals that they desired.

And compared to self-learning in the same environment, the learning of these goals is accelerated, and the resulting representation of these tasks is more useful at a later time. The task sets resulting from guidance are smaller and more closely related to the specific tasks that the human was trying to teach. In self-learning on the other hand, the robot learned a broader task set, serendipitously learning aspects of the environment that the human was not focused on teaching. While not what the human had in mind today, this knowledge about the environment could be advantageous in the future. Clearly both types of learning are beneficial to a robot learner in different ways.

In addition to illustrating the differences between guided and self learning, our experiment allows us to further explore a question that we have raised in prior work: “How do people naturally approach the task of teaching a machine?”

Our video analysis of the learning sessions lets us characterize key similarities and differences in how people use social scaffolding to help teach a physically embodied learner. First, we found that people exhibit consistent behavior in the relative frequency with which they use the different types of scaffolding mechanisms available. Additionally, we were able to learn something about the typical context for each of the scaffolding mechanisms. When making action suggestions or asking Leo to try a known task, people generally wait for eye contact. Presumably waiting for a signal that the robot is finished with its current action and ready to move on. Labeling a state (e.g., “Look, it’s Blue”) mostly happens after eye contact as well, but also happens during an action. Thus, sometimes people want to give the label as soon as the state change happens. Feedback has an even greater split between the eye contact and action contexts. Either a feedback utterance is given at the same time as an action is happening or the person waits until after the action completes and the robot looks up. This raises an important question for a social learning agent. Does a state label or a feedback utterance take on a different connotation when it is given in a different context? It is possible that a person means something different by an utterance given during an action versus one given between actions.

There is some anecdotal evidence that people have different interpretations of how task labeling should work. The current system assumes that the human might provide a label, and if so it would pertain to the current state. Most participants did label in this way, but at least one participant gave ‘pre-labels’ for a given task. Saying, “Leo, now let’s make it Blue.” This is an interaction that the system is currently not designed to handle. Other people gave multiple labels for a state (“Look Leo, it’s open, and it’s green, and the switch is to the right...”). Over half of the participants did this multiple labeling behavior at least once. Again, the system is not designed to take advantage of this, but it is an interesting area for future work. These multiple labels could help the system more quickly learn to differentiate and generalize when a task is considered achieved.

The findings in this study support our previous data with teachable game characters regarding human feedback behavior. Previously, we studied people’s interactions with a virtual robot game character that learns via interactive reinforcement learning (Thomaz and Breazeal 2008). We found that people had a variety of intentions (guidance, motivation) that they communicated in addition to instrumental positive or negative feedback about the last action performed. We found a positive bias in the feedback an agent gets from a human teacher. Additionally, we showed that this asymmetry has a purpose. People mean qualitatively different things with positive and negative feedback. For instance, positive feedback was used to reinforce behavior but also to motivate the agent. Negative feedback was used to “punish” but also to communicate, “undo and back up to your previous

state.”

In the study presented here we see more evidence of the varied nature of positive and negative feedback from a human partner. The split contexts of the feedback messages are an interesting area for future study. It is likely that the feedback takes on a different meaning dependent on the context. Again, we see a positive bias in people’s feedback with 9 out of 11 people using more positive utterances (and in three of those cases the person only gave positive feedback).

The following is an interesting anecdote that highlights the complexity of feedback from a human partner. During a particular learning session, one teacher made the action suggestion to move the switch left when the box switch was already in the left position. Leo did the suggested action, but since it was already left the action had no effect. The teacher gave positive feedback anyway, and then quickly corrected herself and suggested switch right. The true meaning of this positive feedback message is, “yes, you did what I asked, good job, but what I told you was wrong...” Thus, positive and negative feedback from a human partner is much more nuanced than a simple good/bad signal from the environment, and an embodied social learning agent will need the ability to discern these subtle meanings.

A final topic addressed in this study is the extent to which the behavior of the robot influences the human teacher. In prior work, we showed that a virtual robot game character can influence the input from a human teacher with a simple gazing behavior (Thomaz and Breazeal 2006). An embodied robotic agent like Leonardo has many more subtle ways in which to communicate its internal state to the human partner, and we see some evidence that people’s behavior is influenced by the social cues of the robot. On average about 25% of the robot’s facial expressions were mirrored by the human either with their own facial expression, tone of voice, or with a feedback utterance. Also, people waited for Leonardo to make eye contact with them before they would say the next utterance. This has the nice property of a subtle cue the robot uses to slow down the human’s input until the robot is ready for it. In the future, one could imagine exploiting mutual gaze to elicit additional input “just in time.” For instance, the robot might initiate an action, pause and look to the human if confidence is low, to elicit a confirmation or additional guidance before it executes the action.

We see additional anecdotal evidence of people shifting their teaching strategies based on the behavior of the robot. In one case, the person misunderstood the instructions and initially tried to demonstrate the task instead of guide an exploration. She would label a state, describe her actions, and then label the new state. But she quickly shifted into the guided exploration (after about four actions) once Leo started doing actions itself. In another case, the teacher’s strategy was to ‘pre-label’ a task. She would say, “Leo’s let’s make it Blue”, and then make the necessary action suggestions. Once Leo got to the desired state she’d say, “Good Job!” But she did not say the name of the state once they got there, so the label never got attached to that task representation. Then she would ask Leo to make it blue, and Leo would not know the name of that task. Finally, she did one post-labeling, saying the name of the task “Blue” after it was completed, and Leo demonstrated that he could do the blue task soon afterwards. At this point she stopped pre-labeling, and only did the post-labeling for the rest of the learning session.

Thus we see a number of ways that the human partner, in a socially guided learning scenario, actively monitors the robot to build a mental model of the learning process and the effect they are having on it. They adjust their behavior as their mental model solidifies. This supports our approach of having the robot use transparency behaviors to actively communicate with the human partner about

the internal state of the learning process. Additionally, it suggests we continue to develop better mechanisms for helping the human teacher form an appropriate mental model of the robot learner. The robot can use these mechanisms to improve its own learning environment, by improving the quality of instruction received from the human partner.

7. Conclusion

This work acknowledges that a robot learning in a social environment needs the ability to both learn on its own and to take advantage of the social structure provided by a human partner. Our Socially Guided Exploration learning mechanism has motivations to explore its environment and is able to create goal-oriented task representations of novel events. Additionally this process can be influenced by a human partner through attention direction, action suggestion, labeling goal states, and feedback using natural social cues. From our experiments, we found beneficial properties of the balance between intrinsically motivated learning and socially guided learning. Namely, self-exploration tended to result in a broader task repertoire from serendipitous learning opportunities. This broad task set can help to scaffold future learning with a human teacher. Guided-exploration with a human teacher tended to be more goal-driven, resulting in fewer tasks that were learned faster and generalized better to new starting states.

Our analysis of human teaching behavior revealed some interesting findings. First, we found that there was surprisingly little variance among human subjects with respect to how often they used specific types of scaffolding (action suggestions being the highest, negative feedback was the least). Our video analysis reveals different forms of behavior coupling between human teacher and robot learner through social cues. We found that most scaffolding was given to the robot after it made eye contact with the teacher. We also found that human teachers tended to mirror the expressive behavior of the robot (an average of 25%), but this varied by teaching style (some did not mirror at all, some mirrored more than 40%). In addition, we found that the communicative intent behind positive and negative feedback is subtle and varied—it is used in different contexts, sometimes before the robot takes action. Hence, it is not simply reinforcement of past actions. We also found that different teachers have different styles in how they use feedback—some have a positive bias, others are more balanced. Interestingly, none of our subjects had a negative bias.

These findings inform and motivate continued work in how to design robots that learn from human teachers with respect to the dynamic social coupling of teacher and learner to coordinate and improve the teaching/learning process, designing to support the frequency and kinds of scaffolding, understanding the subtlety of intention behind positive/negative feedback, and accommodating different teaching styles.

Acknowledgements

The work presented in this article is a part of ongoing work of the graduate and undergraduate students in the Personal Robotics Group of the MIT Media Lab. The Leonardo robot was initially designed in collaboration with Stan Winston Studios and has been under development since 2002. This work is funded by the Digital Life and Things That Think consortia of the Media Lab, and in particular by the Toyota Motor Corporation.

References

- Chernova, S., and Veloso, M. (2007), "Confidence-Based Policy Learning from Demonstration Using Gaussian Mixture Models," in *Proc. of Autonomous Agents and Multi-Agent Systems (AAMAS)*.
- Atkeson, C.G., and Schaal, S. (1997), "Robot learning from demonstration," in *Proc. 14th International Conference on Machine Learning*, Morgan Kaufmann, pp. 12–20.
- Kaiser, M., and Dillman, R. (1996), "Building Elementary Robot Skills from Human Demonstration," in *Proceedings of the IEEE International Conference on Robotics and Automation*, pp. 2700–2705.
- Kuniyoshi, Y., Inaba, M., and Inoue, H. (1994), "Learning By Watching: Extracting Reusable Task Knowledge From Visual Observation Of Human Performance," *IEEE Transactions on Robotics and Automation*, 10, 799–822.
- Levas, A., and Selfridge, M. (1984), "A User-Friendly High-level Robot Teaching System," in *Proceedings of the IEEE International Conference on Robotics, Atlanta, Georgia*, pp. 413–416.
- Calinon, S., Guenter, F., and Billard, A. (2007), "On Learning, Representing and Generalizing a Task in a Humanoid Robot," *IEEE Transactions on Systems, Man and Cybernetics, Part B. Special issue on robot learning by observation, demonstration and imitation*, 37(2), 286–298.
- Nicolescu, M.N., and Mataric, M.J. (2003), "Natural Methods for Robot Task Learning: Instructive Demonstrations, Generalization and Practice," in *Proceedings of the 2nd Intl. Conf. AAMAS*, July, Melbourne, Australia.
- Lockerd-Thomaz, A., and Breazeal, C. (2004), "Tutelage and Socially Guided Robot Learning," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Sendai, Japan.
- Blumberg, B., Downie, M., Ivanov, Y., Berlin, M., Johnson, M., and Tomlinson, B. (2002), "Integrated learning for interactive synthetic characters," in *Proceedings of the ACM SIGGRAPH*.
- Kaplan, F., Oudeyer, P.Y., Kubinyi, E., and Miklosi, A. (2002), "Robotic clicker training," *Robotics and Autonomous Systems*, 38(3–4), 197–206.
- Saksida, L.M., Raymond, S.M., and Touretzky, D.S. (1998), "Shaping Robot Behavior Using Principles from Instrumental Conditioning," *Robotics and Autonomous Systems*, 22(3/4), 231.
- Clouse, J., and Utgoff, P. (1992), "A teaching method for reinforcement learning," in *Proc. of the Ninth International Conf. on Machine Learning (ICML)*, Aberdeen, Scotland, pp. 92–101.
- Maclin, R., Shavlik, J., Torrey, L., Walker, T., and Wild, E. (2005), "Giving Advice about Preferred Actions to Reinforcement Learners Via Knowledge-Based Kernel Regression," in *Proceedings of the Twentieth National Conference on Artificial Intelligence (AAAI)*, July, Pittsburgh, PA.
- Smart, W.D., and Kaelbling, L.P. (2002), "Effective reinforcement learning for mobile robots," in *In Proceedings of the IEEE International Conference on Robotics and Automation*, Piscataway, NJ, pp. 3404–3410.
- Singh, S., Barto, A.G., and Chentanez, N. (2005), "Intrinsically Motivated Reinforcement Learning," in *Proceedings of Advances in Neural Information Processing Systems 17 (NIPS)*.
- Oudeyer, P.Y., and Kaplan, F. (2004), "Intelligent adaptive curiosity: a source of self-development," in *Proceedings of the 4th International Workshop on Epigenetic Robotics*, Vol. 117, pp. 127–130.
- Schmidhuber, J. (2005), "Self-Motivated Development Through Rewards for Predictor Errors/Improvements," in *Proc. Developmental Robotics 2005 AAAI Spring Symposium*, eds. D. Blank and L. Meeden.
- Lamere, P., Kwok, P., Walker, W., Gouvea, E., Singh, R., Raj, B., and Wolf, P. (2003), "Design of the CMU Sphinx-4 decoder," in *8th European Conf. on Speech Communication and Technology (EUROSPEECH 2003)*, Geneva, Switzerland.
- L. S. Vygotsky, E.M.C., *Mind in society: the development of higher psychological processes*, Cambridge, MA: Harvard University Press (1978).
- Breazeal, C., Berlin, M., Brooks, A., Gray, J., and Thomaz, A.L. (2005), "Using Perspective Taking to Learn from Ambiguous Demonstrations," *Journal of Robotics and Autonomous Systems Special Issue on Robot Programming by Demonstration*.
- Sutton, R., Precup, D., and Singh, S. (1999), "Between MDPs and semi-MDPs: Learning, planning and representing knowledge at multiple temporal scales.," *Journal of Artificial Intelligence Research*, 1, 139.
- Breazeal, C., *Designing Sociable Robots*, Cambridge, MA: MIT Press (2002).
- Sutton, R., Precup, D., and Singh, S. (1998), "Intra-option learning about temporally abstract actions.," in *Proceedings of the Fifteenth International Conference on Machine Learning (ICML98)*, Masion, WI.
- Clark, H.H., *Using Language*, Cambridge: Cambridge University Press (1996).
- Smith, C., and Scott, H. (1997), "A Componential Approach to the meaning of facial expressions," in *The Psychology of Facial Expression* United Kingdom: Cambridge University Press.
- Thomaz, A.L., Berlin, M., and Breazeal, C. (2005), "An Embodied Computational Model of Social Referencing," in *IEEE International Workshop on Human Robot Interaction (RO-MAN)*, Nashville, TN.
- Thomaz, A.L., and Breazeal, C. (2008), "Teachable Robots: Understanding Human Teaching Behavior to Build More Effective Robot Learners.," *Artificial Intelligence Journal*, 172, 716–737.
- Thomaz, A.L., and Breazeal, C. (2006), "Transparency and Socially Guided Machine Learning," in *Proceedings of the 5th International Conference on Developmental Learning (ICDL)*, Bloomington, IN.