

# Goals vs. Rewards: A Comparative Study of Objective Specification Mechanisms

**Anonymous authors**

Paper under double-blind review

**Keywords:** objective specification, goals, rewards.

## Summary

This paper looks at two popular objective specification mechanisms for sequential decision-making problems: goals and rewards. We investigate how easy it is for non-AI experts to use these different specification mechanisms effectively. Specifically, through this paper, we investigate how effectively these mechanisms could be used to (a) correctly direct an AI system or robot to generate some desired behavior and (b) predict the behavior encoded in a given objective specification. We perform a user study to assess these questions. In addition, we present a formalization of the problems of objective specification and behavior prediction, and we characterize underspecification and overspecification. While participants have a strong preference for using goals as an objective specification mechanism, we find a surprising result: even naïve users are equally capable of specifying and interpreting reward functions.

## Contribution(s)

1. The paper compares and contrasts how well naïve users can effectively make use of goal and reward specification mechanisms. In particular, we assess whether they (a) can use these mechanisms to generate specifications that can result in some intended target behavior and (b) whether they can predict behavior that could result from the given specification.  
**Context:** We are unaware of any works that perform such human-centric comparisons. The closest works we know of focus purely on how successful engineers are in hand-crafting reward functions (cf. (Knox et al., 2023; Booth et al., 2023)).
2. We provide a formal definition of the specification and prediction task to support comparisons between reward functions and goals. We also provide a formal characterization of the conditions under which an objective can be said to be overspecified or underspecified.  
**Context:** While there are existing works that have tried to model objective misspecification (e.g., Mechergui & Sreedharan (2024)), underspecification (e.g., Shah et al. (2022)), and misspecification (e.g., Amodei et al. (2016)), these definitions have not been formalized to cover and compare multiple specification modalities.
3. Our results present evidence that the naïve users’ ability to correctly specify and interpret reward functions is comparable to their ability to provide goal specifications. However, we see a clear difference in their preferences between the two metrics: they overwhelmingly prefer the goal mechanism.  
**Context:** We are unaware of any prior works that point to parity in user ability to leverage the two objective specification mechanisms. This result may imply that developing novel interfaces for reward functions could help users of RL techniques to utilize reward functions more effectively—for example, reward machines are one such promising mechanism (Icarte et al., 2022).

# Goals vs. Rewards: A Comparative Study of Objective Specification Mechanisms

**Anonymous authors**

Paper under double-blind review

## Abstract

1 This paper looks at two popular objective specification mechanisms for sequential  
 2 decision-making problems: goals and rewards. We investigate how easy it is for peo-  
 3 ple without AI expertise to use these different specification mechanisms effectively.  
 4 Specifically, through this paper, we investigate how effectively these mechanisms could  
 5 be used to (a) correctly direct an AI system or robot to generate some desired behav-  
 6 ior and (b) predict the behavior encoded in a given objective specification. We first  
 7 present a formalization of the problems of objective specification and behavior predic-  
 8 tion, and we characterize underspecification and overspecification. We then perform a  
 9 user study to assess how well participants are able to use rewards and goals as speci-  
 10 fication mechanisms, and their propensity for overspecification and underspecification  
 11 with these mechanisms. While participants have a strong preference for using goals as  
 12 an objective specification mechanism, we find a surprising result: even naïve users are  
 13 equally capable of specifying and interpreting reward functions as of using goals.

## 14 1 Introduction

15 We examine the two common specification mechanisms for sequential decision-making: goals and  
 16 rewards. We assess how well non-AI experts can work with these different specification mecha-  
 17 nisms. Goals and rewards have different expected upsides. Goals allow people to provide a partial  
 18 specification of their desired end states. This mechanism is commonly used in classical planning  
 19 (Cox, 2016) and has also received a lot of attention from recent work in using Large Language Mod-  
 20 els (LLMs) (Brown et al., 2020) for robot planning (cf. (Brohan et al., 2023)). Rewards, on the  
 21 other hand, are the underlying objective specification mechanisms used by reinforcement learning  
 22 (RL) methods (Sutton & Barto, 2018) and Markov Decision-making Processes (MDPs) Rewards are  
 23 a means for encoding goals: the reward hypothesis asserts that “all of what we mean by goals and  
 24 purposes can be well thought of as maximization of the expected value of the cumulative sum of a  
 25 received scalar signal [reward]” (Sutton & Barto, 2018). The reward format allows one to associate  
 26 scalar signal with reaching some state or performing some action in a given state.

27 The research community has developed a rigorous understanding of these specification mechanisms’  
 28 expressiveness and representational limitations (cf. (Abel et al., 2021)). Despite this understanding,  
 29 the ease with which users can express their underlying objectives in these forms has not, to our  
 30 knowledge, been explicitly studied. While the development of LLMs has received attention as po-  
 31 tentially intuitive interfaces to AI systems, they do not entirely address the question of how to best  
 32 construct specifications for AI systems either. After all, LLMs would need to translate the user  
 33 utterances into the underlying objective specification (whether goals, rewards, or some other speci-  
 34 fication form), and it is unclear if these utterances would contain sufficient information needed for  
 35 the translation.

36 In this paper, we conduct a user study to examine the ease of use, strengths, and weaknesses of  
 37 the two specification mechanisms when used by non-AI experts. In the user study, we expose par-

Participants to these objective specification mechanisms in intuitive tasks using simple interfaces and measure (a) how well the users are able to use the specific mechanism correctly and (b) how well they can understand an objective specified using each mechanism. While there have been some efforts at measuring the difficulty in specifying rewards (Booth et al., 2023), to the best of our knowledge, our work represents the first effort to perform such a comparative analysis of the two specification mechanisms among non-AI experts. To ensure our user studies are performed from a firm formal grounding, we also provide a concrete characterization of the tasks related to objective specification and behavior prediction given an objective. Additionally, we provide a characterization for when a given objective could be said to be over or under-specified.

The two primary takeaways from our study results are as follows. First, naïve users are not as bad at specifying reward functions as is generally assumed, and in fact their ability to do so is comparable to their ability to correctly specify goals. This is a surprising result, as goals are a seemingly more intuitive mechanism and are more commonly represented in everyday communications. Second, despite their ability to use rewards as specifications, users generally perceive goal specification to be more intuitive and easier to specify. We believe that the results from this study could help us design objective specification interfaces that are more intuitive and easy to use for everyday users.

The paper is structured as follows: Section 2 discusses the related works. Section 3 provides a brief discussion of goals and rewards as an objective specification mechanism and potential trade-offs. Section 4 describes the formal definition of specification and prediction. We describe the specific hypotheses we focus on in Section 5. Section 6 discusses the methods, including the study design. Section 7 presents the results and discussions. Finally, the conclusion is described in Section 8.

## 2 Related Work

The notion that goals are a natural way people think about their objectives has a long history. One could see similar ideas being discussed Aristotle’s notions of *phronēsis* (Taylor, 2019) to means-end analysis (Simon, 2019). Apart from evidence that people may leverage some notions of goals in their own reasoning, there have been fewer studies performed in determining if goals are, in fact, the best mechanisms for people to actually specify their objectives. Some works within this space include proposals to compare how effectively people can specify their objectives in procedural terms, i.e., in terms of actions or sequence of actions, as opposed to the end goal (Tran, 2024).

In the reward space, reinforcement learning often assumes the existence of a divined reward function that encodes the task. In practice, though, correctly specifying reward functions is nontrivial: the challenge of doing so correctly has catalyzed the take-off of the AI safety research community (Amodei et al., 2016; Russell, 2022). Further, reward functions are typically designed by engineers through trial-and-error design processes (Knox et al., 2023), which are subject to oversights and inaccuracies, even when crafted by reinforcement learning experts (Booth et al., 2023).

Because of the challenges of using either goals or rewards as specifications, efforts in human-computer interaction, broadly construed, have sought to use intuitive signals in place of these explicit specification modalities. These alternatives span feedback (Knox & Stone, 2009; MacGlashan et al., 2017), corrections (Losey & O’Malley, 2018; Bajcsy et al., 2018), advice (Thomaz & Breazeal, 2008; Amershi et al., 2014), demonstrations (Ravichandar et al., 2020), dynamical system modulation matrices (Figueroa et al., 2020), and, most famously, preferences (Christiano et al., 2017; Ziegler et al., 2019; Biyik & Sadigh, 2018). While these intuitive mechanisms unlock naïve users’ ability to program machines, their interpretation is subject to failures and misinterpretation since the human providing the specification has less control over how the system interprets their specification. For example, a line of research has questioned the inductive bias used in reinforcement learning from human preference (Knox et al., 2022).

### 84 3 Background

85 We will start by providing a brief sketch of the two specification mechanisms under consideration,  
 86 goals and rewards. Since we primarily focus on sequential decision-making settings, for each prob-  
 87 lem, we will separate out the the task domain from the objective specification. In each case, the task  
 88 domain will provide the details on the dynamics of the task and the starting state of the environment.

89 To start with, goals as an objective specification mechanism is most commonly used in deterministic  
 90 factored planning settings, also referred to as “classical planning” settings (Geffner & Bonet, 2013).  
 91 In general, a classical planning problem can be represented by a tuple of the form  $\mathcal{P}^c = \langle \mathcal{D}^c, \mathcal{G}^c \rangle$ ,  
 92 where  $\mathcal{D}^c$  is the task domain and  $\mathcal{G}^c$  is the goal specification. The task domain is further defined as  
 93  $\mathcal{D}^c = \langle F^c, A^c, I^c \rangle$ , where  $F^c$  is a set of proposition variables or facts used to define the state space,  
 94  $A^c$  is the set of actions and  $I^c$  is the initial state. Each action  $a \in A^c$ , is further defined by a tuple of  
 95 the form  $a = \langle pre(a), add(a), del(a) \rangle$ . Here  $pre(a) \subseteq$  is the preconditions that need to be satisfied  
 96 for the action  $a$  to be executable,  $add(a)$  and  $del(a)$  are add and delete effects, respectively. The  
 97 result of executing an action  $a$  in state  $s$ , is captured by the transition function  $\Gamma^c$ , and is given as:

$$\Gamma^c(s, a) = \begin{cases} (s \setminus del(a)) \cup add(a) & \text{If } pre(a) \subseteq s \\ Undefined & \text{Otherwise} \end{cases}$$

98 We will also overload the notation and use  $\Gamma^c$  to denote the execution of action sequences. A solution  
 99 to a classical planning problem takes the form of an action sequence whose execution in the initial  
 100 state results in a state that satisfies the goal specification. Such an action sequence is referred to as a  
 101 plan. More formally, an action sequence  $\pi = \langle a_1, \dots, a_k \rangle$  is a plan if  $\Gamma^c(I^c, \pi) \supseteq \mathcal{G}^c$ . In the simplest  
 102 formalism, an optimal plan corresponds to the shortest possible plan<sup>1</sup>.

103 Reward functions are defined in the context of a Markov Decision Process or MDP (Puterman,  
 104 1990). Here, an MDP will be defined using a tuple of the form  $\mathcal{P}^m = \langle \mathcal{D}^m, \mathcal{R}^m \rangle$ . As with the  
 105 previous planning formalism,  $\mathcal{D}^m$  stands for the domain, but our objective is now given by a reward  
 106 function  $\mathcal{R}^m$ . In this case, the domain is given by a tuple of the form  $\langle F^m, A^m, I^m, T^m, \gamma \rangle$ , now as  
 107 before  $F^m$  stands for the state variable and  $I^m$  the initial state. Here,  $A^m$  only lists the action labels,  
 108 and the dynamics of the action are determined completely by the transition probability function  
 109  $T^m$ . Finally,  $\gamma \in [0, 1)$  represents the discount factor that determines how the agent maximizes  
 110 cumulative discounted future rewards or returns. Here, we will also have a slightly different state  
 111 space. Specifically, we will define it as  $S^m = 2^F \cup \{\perp\}$ . Here, we add the new state  $\perp$  as a stand-in  
 112 for the end state. Now, the transition function will be given as

$$T^m : S^m \times A^m \times S^m \rightarrow \{0, 1\}$$

113 Here, the mapping is only to probabilities 0 and 1 since we focus on problems with deterministic  
 114 transition probabilities. To support the transition into end states, we will also introduce an exit action  
 115  $\mathcal{E} \in A^m$ , that will deterministically transition into the end state  $\perp$ .

116 We will define the reward function as  $\mathcal{R}^m : F \times A \rightarrow \mathbb{R}$ , i.e., a mapping from a state variable and  
 117 action pair to a number. The reward associated with a state, action pair is given as

$$\mathcal{R}^m(s, a) = \begin{cases} \sum_{f \in S} \mathcal{R}^m(f, a) & \text{if } s \neq \perp \\ 0 & \text{otherwise} \end{cases}$$

118 A solution to an MDP problem takes the form of a policy  $\pi : S^m \rightarrow A$ , i.e., a function that maps  
 119 states to actions. A policy is said to be optimal if it maximizes the total expected discounted reward  
 120 received under the given policy.

121 At this point, it is worth noting that for every classical planning task domain  $\mathcal{D}^c$ , we can build a  
 122 corresponding task domain  $\mathcal{D}_m^c = \langle F_m^c, A_m^c, I^c, T_m^c, \gamma \rangle$ , where  $F_m^c = F^c \cup \{\perp\}$ ,  $A_m^c$  one action

<sup>1</sup>However, there are more expressive formalisms that allow one to associate non-unit costs with actions

label for each action in  $A^c$  plus a label for  $\mathcal{E}$ ,  $I^c$  is the initial state (and same as before), the transition  $T_m^c$  returns one only if it is a valid transition per  $\Gamma^c$ . For the application of actions in states where the preconditions aren't met, we will assign a probability of '1' to transition to  $\perp$ , and  $\perp$  is treated as an absorber state.

We will use the notion of trace as a shared notion of behavior that can be used in both settings. A trace  $\tau$  for a policy or plan consists of a sequence of state action pairs that results from the execution of a policy or plan in the initial state. We will also use notation  $\mathcal{P} = \langle \mathcal{D}, \mathcal{O} \rangle$  as a generalized scheme of model representation that can stand in for both classical planning problems and MDP. Depending on the context,  $\mathcal{O}$  could either be a reward or a goal.

## 4 Specification and Prediction

With the basic notations in place, we can precisely define the exact questions under examination. In particular, we are interested in the user's ability to specify an objective that can lead to some desired behavior or be able to predict behavior that could result from optimizing for a given objective function. These two problems correspond to the primary ways users specify objectives. We start with the specification problem, where a user must identify an objective resulting in a target behavior.

**Definition 1** For a given domain model  $\mathcal{D}$  and a target trace  $\tau$ , the specification problem corresponds to finding an objective  $\mathcal{O}$ , such that  $\tau$  is a trace for an optimal solution for the problem  $\mathcal{P} = \langle \mathcal{D}, \mathcal{O} \rangle$ .

If the optimal solution for a given objective specification (i.e., a goal or reward) leads to a trace  $\tau$ , then we will refer to that objective specification as being a correctly specified objective for  $\tau$ , else it is referred to as a misspecified objective.

Moving from the more general to specific settings, we start seeing differences in properties. For example, one can show that even when a goal specification cannot be found for a given trace, it might be possible to find a reward function in a corresponding MDP.

**Proposition 1** For a classical planning domain  $\mathcal{D}^c$ , let  $\tau = \langle I, a_0, \dots, s_k \rangle$ , be a trace such that for every consecutive state-action-state tuple  $s_i, a_i, s_{i+1}$  we have  $\Gamma^c(s_i, a_i) = s_{i+1}$ , and the trace contains no repeating states, then even if there exists no goal for which  $\tau$  is a trace for an optimal plan, there still exists a reward function for  $m(\mathcal{D}^c)$  for which  $\tau$  is a trace for an optimal policy.

The above proposition can be proven by showing that there exist traces that satisfy the property for which no goal exists and by showing the existence of a reward for which the trace is part of an optimal policy. For the first, consider a trace that includes an avoidable subsequence. In other words, let  $s_i$  and  $s_j$  be part of  $\tau$  such that their positions in sequences are separated by more than two positions, i.e., there are at least two actions between  $s_i$  and  $s_j$ . Now let's assume there exists an action  $a$ , such that  $\Gamma^c(s_i, a) = s_j$ . Then, by definition, this trace can't be part of an optimal plan since you can get a shorter trace that results in the same state by removing the original actions between  $s_i$  and  $s_j$ . As for the second part, consider a reward function that assigns zero to every state. Under this reward function, all policy has the same value and are optimal. Given the fact that all transition in the trace corresponds to valid ones in the original domain model, there exists at least one policy for which this is a valid trace.

This clearly shows how the reward function provides a clear advantage in terms of expressivity. However, this advantage goes even further: the knowledge about the goal will allow us to reconstruct a reward function for the corresponding model directly. Specifically, one can create a reward function that assigns a positive reward to all the goal fluents for the exit action, or more formally,

166 **Proposition 2** For a trace  $\tau$  and a classical planning domain  $\mathcal{D}^c$ , let  $\mathcal{G}^c$  be a correctly specified  
 167 goal, then  $\mathcal{R}_m^c$  must be a correctly specified reward for  $m(\mathcal{D}^c)$ , when

$$\mathcal{R}_m^c(f, a) = \begin{cases} r^+ & \text{if } f \in \mathcal{G}^c \text{ and } a = \mathcal{E} \\ 0 & \text{otherwise} \end{cases}$$

168 The validity of the above proposition is straightforward. The agent only receives a positive reward  
 169 for performing exit action from states that satisfy the goal specification. The presence of a discount  
 170 factor means that this would need to be achieved in as few steps as possible.

171 Now, it is also worth noting that not all correctly specified objectives are equal. In particular, we  
 172 can identify two categories. In one case, the user may not have provided enough details; we will  
 173 call such cases examples of underspecification. In the latter case, the user would have provided  
 174 more details than needed or examples of overspecification. It is worth noting that the implications  
 175 of the two are widely different. While overspecification might reduce the set of optimal policies  
 176 and prevent the AI system from coming up with creative solutions, underspecification could result  
 177 in unexpected behavior or specification gaming. We can define the two categories as follows:

178 **Definition 2** For a domain model  $\mathcal{D}$  and a target trace  $\tau$ , a given specification  $\mathcal{O}$  is said to be  
 179 underspecified if there are other traces  $\tau' \neq \tau$  that could result from other optimal solutions for  
 180  $\mathcal{P} = \langle \mathcal{D}, \mathcal{O} \rangle$ .

181 In the above definition, underspecification is purely defined by the fact there are other traces and  
 182 solutions possible (given the deterministic settings we consider, there is a one-to-one mapping be-  
 183 tween solutions and traces). On the other hand, defining overspecification requires us to use a notion  
 184 of specification size, i.e.,  $|\mathcal{O}|$ , where for goals, the size is given by the number of fluent in the speci-  
 185 fication and for rewards, the number of fluent action pairs with non-zero values. Now, we can define  
 186 overspecification to be cases where specifications of smaller size exist that are not underspecified.

187 **Definition 3** For a domain model  $\mathcal{D}$  and a target trace  $\tau$ , a given specification  $\mathcal{O}$  is said to be  
 188 overspecified if (a)  $\mathcal{O}$  is not underspecified and (b) there exists another correct specification  $\mathcal{O}'$ ,  
 189 such that  $\mathcal{O}'$  is not an under specification and  $|\mathcal{O}'| < |\mathcal{O}|$ .

190 This brings us an end to the section discussing the first task, namely objective specification. The sec-  
 191 ond task corresponds to the user's ability to make inferences based on the given objective. Here, we  
 192 consider the simple case of whether a user can tell if a trace is possible under a given specification.

193 **Definition 4** For a given problem  $\mathcal{P} = \langle \mathcal{D}, \mathcal{O} \rangle$  and a trace  $\tau$ , the prediction problem corresponds  
 194 to identifying whether  $\tau$  is a trace for an optimal solution for the problem  $\mathcal{P}$ .

## 195 5 Hypotheses

196 Our study is primarily designed to measure how the choice of specification mechanism can affect  
 197 the user's ability to specify objectives and predict agent behavior. The primary hypotheses we plan  
 198 to test here are as follows:

- 199 • H1-a: Participants are more likely to provide accurate goals than accurate reward specifications.
- 200 • H1-b: Participants are more likely to correctly interpret goals than reward specifications.

201 The next question we consider concerns the participants' workload—in particular, the cognitive  
 202 load, imposed and the time taken by the two mechanisms.

- 203 • H2-a: Reward specifications will result in a higher workload than goal specifications and will  
 204 require longer time to finish.
- 205 • H2-b: Trying to interpret reward functions will result in a higher workload than goal specifications  
 206 and will require a longer time to finish.

Now, we also wanted to use this as an opportunity to understand ways in which the user specification may differ from the minimal specification, which brings us to the hypothesis:

- H3: Participants are more likely to underspecify objectives than overspecify them.

We will test the above hypothesis for both reward and goal specification cases.

To assess the H2, we measure the participants' workload for each objective specification mechanism and task in the survey. NASA Task Load Index (TLX) is used to measure the perceived workload. NASA TLX has six dimensions: mental demand, physical demand, temporal demand, performance, effort, and frustration level (Hart, 1986). Each dimension is measured using a Likert rating scale.

## 6 Methods

### 6.1 Study Design

To compare the two mechanisms, we designed three intuitive but diverse domains in which two primary tasks related to each mechanism can be tested: (1) the user's ability to provide an objective specification that will result in a given behavior and (2) their ability to predict the behavior from a given specification. We chose domains that non-AI experts could understand without considerable training but corresponded to potential real-world robotics applications. Specifically, the domains included (1) a robot navigation task, (2) a tabletop pick-and-place task, and (3) a task with a self-driving vehicle. We chose deterministic versions of the tasks to avoid potential confounders that may arise from the stochasticity of the environment dynamics. The environment setting for each domain can be seen in Figure 1.

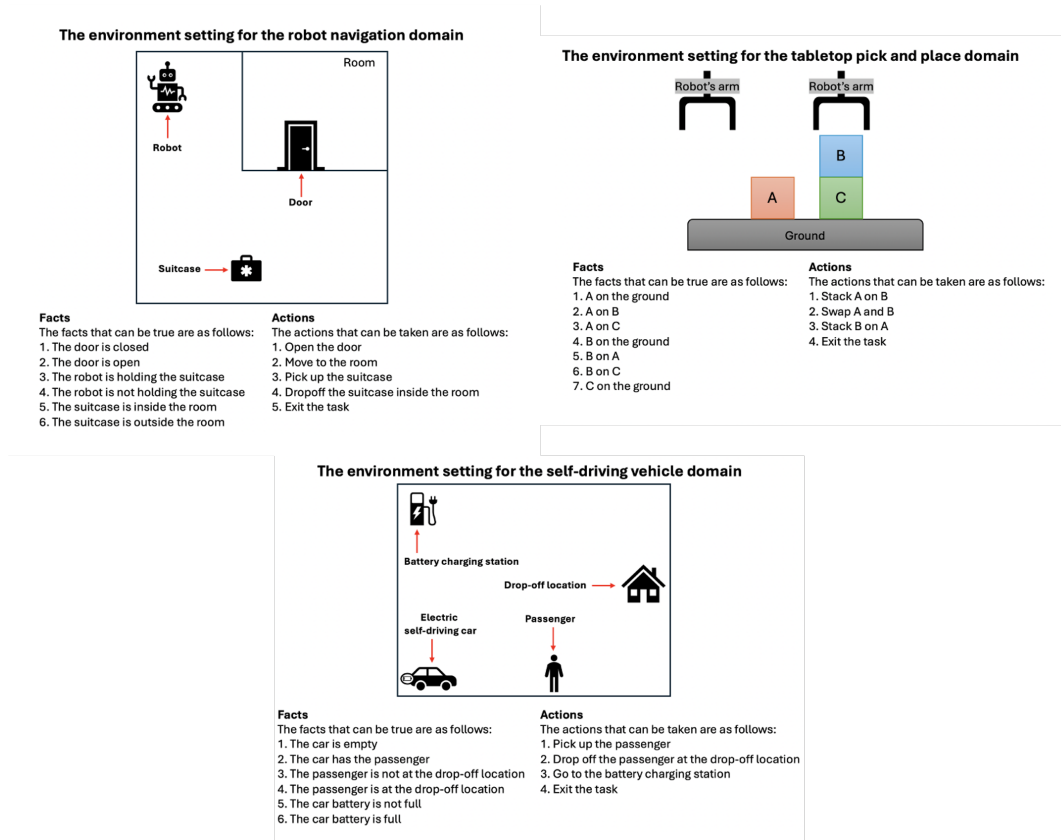


Figure 1: A visualization of each domain used in the study. Top left: a robot navigation task. Top right: a pick-and-place task. Bottom: a self-driving vehicle task.



The navigation task involves robots navigating through a workspace. In this case, a robot needs to pick up and drop off a suitcase in different locations within a small workspace. The pick and place domain contains a set of blocks that can be stacked on top of one another. The objective is usually to achieve a specific configuration of the blocks. For the self-driving vehicle domain, we have a self-driving car powered by a battery that needs to pick up and drop off a passenger in different locations. It also needs to charge the battery to make sure that the battery is enough to perform its task. In each environment setting, the current state is defined by a set of binary variables, henceforth referred to as facts. There is also a set of actions that can be taken by the robot, including an exit action that will allow the robot to end the task. Each domain had about 6-7 facts and 4-5 actions. We choose to keep the facts and action counts similar so as to balance the workload between domains.

We use these domains to create surveys that test the participants' ability to specify an objective that will result in some provided behavior or their ability to predict what behavior will result from a given objective. The survey built around these scenarios uses a mixed study design, combining both between-subjects and within-subjects study designs. The participants are shown either the specification task or prediction task (making this study design between subjects), chosen from three different problem domains as mentioned above. Given the problem domain, the participants are tested on how well they are able to complete the specified task across the two objective specification mechanisms (within subjects). We will use a counterbalancing technique to vary the order in which participants will be shown the different specification mechanisms. This is to ensure that no single order influences the results of the study.

For each objective specification mechanism, there are two sections in the survey: demo and test. The demo section is basically a learning phase, where participants are familiarized and introduced to the concepts of goal and reward specifications. In the demo section, participants will be shown a video that demonstrates a simple behavior along with the corresponding goal or reward (see the example illustration in Figure 2). For goals, the video will show the "facts to be achieved (goal state)" and how the "facts that are true (current state)" change during the duration of robot behavior until it reaches the goal state. On the other hand, for rewards, the video shows the rewards matrix and how individual rewards from the matrix will be added to the total when the agent performs specific actions. For example, based on the illustration in Figure 2, the agent will get 50 points if it takes an "exit the task" action while the fact that "the robot is holding the suitcase" is true.

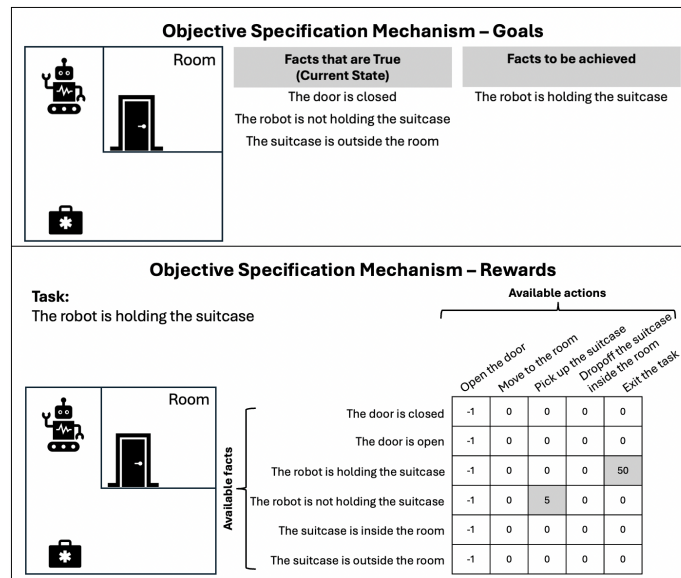


Figure 2: Illustrations for the sample specifications that could be shown to the participants.



For the first task, i.e., ease of objective specification, the test section will show a sample behavior to the user. Then, participants are asked to come up with goals and/or rewards for that scenario. Figure 3 presents screenshots of the interface provided to the user to specify the objective. We refer to goals as facts and rewards as scores to simplify the description to non-AI expert participants. From the participants’ answers, we can determine whether their specifications are correct or incorrect. For the incorrect one, the potential sources of errors can be analyzed, including overspecification and underspecification.

**Objective Specification Mechanism - Facts**

Drag the adequate Fact(s) that will allow us to recreate the same behavior seen in the previous video and drop them into the Facts to be achieved. Try to list the minimal number of facts to be achieved, that will generate the expected behavior.

Items

- The door is closed
- The door is open
- The robot is holding the suitcase
- The robot is not holding the suitcase
- The suitcase is inside the room
- The suitcase is outside the room

Facts to be achieved

**Objective Specification Mechanism - Scores**

Each cell in the table of scores will have a default value of zero. For the specification, update the relevant cells. If they are not updated, they will contribute zero points to the total score. The range of values that you can put in for each cell is a number between -100 and 100. Try to fill in the minimal number of cells that will generate the expected behavior.

	Open the door	Pick up the suitcase outside the room	Dropoff the suitcase inside the room	Exit the task
The door is closed	0	0	0	0
The door is open	0	0	0	0
The robot is holding the suitcase	0	0	0	0
The robot is not holding the suitcase	0	0	0	0
The suitcase is inside the room	0	0	0	0
The suitcase is outside the room	0	0	0	0

Figure 3: Sample interfaces used by the participants to specify goals and rewards. The one shown above corresponds to the navigation task.

On the other hand, to test how easily non-AI experts can understand goals and rewards, instead of showing the demonstration, we show the correct goal (list of facts to be achieved) or the rewards specification (in the form of scores). Then, we ask the participants to predict or interpret the behavior of the agent based on that. Specifically, we provide three video options and ask them to choose one that most aligns with the given goals or rewards.

Additionally, at the end of the survey, we ask the participants to directly compare the two specification mechanisms in terms of their easiness, intuitiveness, likeability, and challenge. We also ask for qualitative feedback on why they think that particular objective specification mechanism is easier or harder than the other. Finally, we collect demographic information such as age, gender, highest level of education, and familiarity with computer science and AI subjects.

## 6.2 Participants and Procedure

Before the main study, we ran a small pilot think-aloud study with three participants to refine the study design. For the primary user study, we recruited a total of 30 participants from Prolific: 15 participants (8 males and 7 females) for the specification task and 15 participants (7 males and 8 females) for the prediction task. They were paid \$18.5 USD per hour, and they identified their native language as English. The majority of them reported having never taken an AI course.

This study was IRB-approved. Participants were provided with informed consent before they started the survey. Multiple attention check questions were included throughout the study. Each participant was shown all of the three domains in random order. The order in which the specification mechanism was shown was also randomized to ensure the results were counterbalanced.

## 7 Results and Discussions

### 7.1 Impressions from the Think-Aloud Study

We used the think-aloud study (Baxter et al., 2015) as a means of both testing our interface, particularly for specification tasks and collecting some initial anecdotal information on the mechanisms. The reactions we observed were aligned with what we hypothesized (H1-a and H2-a), where the participants showed more positive reactions to the goal specification interface as opposed to the re-

ward. Some reactions to goal specifications included: “This one is fun, like playing,” and “The task was super easy.” On the other hand, for the reward specification, users reported a lack of confidence about their ability to correctly provide such specifications: “I don’t understand, I’m very bad at this,” and “I don’t know why this is confusing me.” Their qualitative feedback at the end of the survey also reflected their strong preference for using the goal specification mechanism.

## 7.2 Specification Task

We started by analyzing the initial results from the specification task. In regards to hypothesis H1-a, we calculated the number of times the participants were able to provide correct specifications (presented in Table 1). We were surprised to find that the participants were actually able to identify correct reward functions more frequently than correct goal specifications. Further, analyses of the results showed that the most frequent mistakes made by participants in goal specification involved the inclusion of intermediate facts in the goal specification. These intermediate facts, while made true by the agent’s action, are also made false by further actions in the plan. For example, the subjects might indicate that “the robot is holding the suitcase”—but, in the observation of the environment, the robot places the suitcase down at the end of the video. As such, including these intermediate facts in the final goal specification leads to an unachievable objective specification. The goals being provided by the users reflected a more procedural description of the agent behavior than a final goal state description. On the other hand, such intermediate state scores can be more naturally incorporated into the reward function. To analyze the factors that could explain these results, we created two follow-up variants of the specification and reran our study.

In the first follow-up, we updated all our videos to highlight how intermediate fact values change. In each of our demonstration videos, we added animations that showed which facts became false. We reran the experiment on five participants (thus collecting 15 specifications per mechanism). The results from the study are presented in Table 2. While we see that the additional information does improve the overall percentage of correct goal specification, the resulting percentage is similar to that of the rewards—indicating that this additional information balances participants’ ability to craft goals or rewards.

In the second follow-up, we considered a variation of the navigation task where the participants simply provided scores for each state variable achievement. This variant was motivated by the possibility that the inclusion of actions in the specification mechanism might be helping the participant by allowing them to think procedurally about the task. Here, we set a specific absorbing state, and the reward for each state was set to the sum of rewards associated with each state factor. We ran this variant on 15 participants, and the percentage of correct goal and reward specifications were, in fact, the same (Table 3). This shows that the presence of actions assisted participants in crafting rewards, and that crafting rewards over states instead of state-action-state tuples is a harder task in these domains.

Taken together, this collection of results points to the possibility that the hypothesis that goals are easier than rewards to specify need not be true. This is particularly surprising, given this hypothesis is quite frequently taken to be self-evident in the literature (cf. (Mechergui & Sreedharan, 2024)).

However, when we move on to the hypothesis related to workload and time taken (H2-a), we see a clear distinction between the two specification mechanisms, with subjects overwhelmingly preferring the goal mechanism. Running paired t-tests shows that there is a statistically significant difference between the cognitive load of goal specification ( $M = 9.444$ ,  $SD = 5.057$ ) and reward specification ( $M = 12.689$ ,  $SD = 5.008$ ). There is also a statistically significant difference between the time taken to complete the goal specification ( $M = 82.014$ ,  $SD = 36.225$ ) and reward specification ( $M = 148.521$ ,  $SD = 87.015$ ). In addition, we also see similar responses with respect to the qualitative responses, with most participants finding goals easier to specify (86.67%) and more intuitive (73.33%). The supplementary file provides the breakdown of individual dimensions of the workload and more details on the qualitative feedback. These results support our hypothesis H2-a.

Finally, moving to H3, our results again do not support our hypothesis. In fact, we saw more instances of the users overspecifying their objectives than underspecification (see Table 1, 2, and 3). Such patterns were also replicated in the incorrect specifications. Looking at incorrect goal specification, we saw a larger set of participants (75%) added incorrect facts as opposed to leaving out some facts (0.027%).

Table 1: Results from the main specification user study

Category	Sub-category	Percentage of total response	
		Goals	Rewards
Correct	Correct minimal specification	4.45	2.22
	Correct but overspecified	13.33	35.56
	Correct but underspecified	-	8.89
Incorrect	Incorrect because gave subset	82.22	53.33
Total		100	100

Table 2: Results from the first variant of the specification study

Category	Sub-category	Percentage of total response	
		Goals	Rewards
Correct	Correct minimal specification	-	-
	Correct but overspecified	73.33	66.67
	Correct but underspecified	-	-
Incorrect	Incorrect specification	26.67	33.33
Total		100	100

Table 3: Results from the second variant of the specification study

Category	Sub-category	Percentage of total response	
		Goals	Rewards
Correct	Correct and minimal specification	-	-
	Correct and overspecification	66.67	66.67
	Correct and underspecification	-	-
Incorrect	Incorrect specification	33.33	33.33
Total		100	100

### 7.3 Prediction Task

As discussed, the goal of the prediction task was to test whether a user can predict the behavior that could result from a given specification. We see that as a proxy for the ease with which users can correctly interpret specifications expressed using each mechanism. For the prediction task, we also see a similar pattern. The participant accuracy in predicting behavior based on the given goals function and reward function is comparably high, with 93.33% predicting the goals function correctly and 91.11% predicting the rewards function correctly. Here, the difference is not high enough to establish any statistically significant difference between the two groups. As for the results related to the cognitive workload, our t-test was not able to establish any significant difference between the prediction from the goal function ( $M = 6.267$ ,  $SD = 6.308$ ) and from the reward function ( $M = 7.044$ ,  $SD = 6.502$ ) with  $P$ -value equal to .251. There was also no significant difference

354 between the time taken to complete the prediction from goal function ( $M = 82.840$ ,  $SD = 75.285$ )  
355 and from reward function ( $M = 90.873$ ,  $SD = 59.339$ );  $t(44) = -0.878$ ,  $P = .385$ . This seems to  
356 suggest that both of our hypotheses H1-b and H2-b may not hold.

357 However, when we move on to the participant preference between the two mechanisms, most par-  
358 ticipants find goal function easier to predict (86.67%) and more intuitive (80%). In addition to this,  
359 most participants reported that the reward function is more challenging to predict (80%). These  
360 preferences are consistent in both specification and prediction tasks.

## 361 7.4 Results Summary

362 From our experiments, we find that our hypotheses H1-a and H1-b are surprisingly not supported:  
363 we did not find evidence that people are able to more correctly specify or interpret goals over reward  
364 functions. Despite this, we find that there was a significant difference in the cognitive effort and  
365 time needed to specify objectives: goals were a clear winner on these axes (H2-a). We were also  
366 surprised to find that people are not more likely to underspecify goals than to overspecify (H3).  
367 Overall, though, the subjective feedback reflects that participants strongly preferred using goals  
368 over reward functions.

## 369 7.5 Limitations of Study Scenarios

370 Please note that all studies were carried out in purely deterministic settings, where the agents can  
371 not get stuck in loops. While this is stereotypical of many tasks where goals are used, this doesn't  
372 necessarily represent all the ways rewards could be utilized, which is a more general specification  
373 mechanism. Similarly, we considered simple enough scenarios where the participants could easily  
374 enumerate all possible facts and incorporate them into the specification.

## 375 8 Conclusion

376 In this paper, we performed a comparison to assess how easy it would be for naïve users to provide  
377 and understand reward specifications and goal specifications. Our results point to the fact that,  
378 in fact, people's ability to provide and understand rewards is fairly comparable to that of goals.  
379 However, there is a clear difference in the user preferences and the cognitive load imposed by the  
380 two methods (at least for the specification task). One interesting question to ask in this context would  
381 be whether this difference can be explained by the interface we used for our study. As such, one  
382 would want to investigate if it's possible to develop interfaces that allow users to intuitively provide  
383 reward functions. Such interfaces would have pretty immediate advantages, given reward functions  
384 are more expressive than goals. Also, as the next steps, we would also like to investigate how  
385 goals and rewards compare against other objective specification mechanisms like policy sketches  
386 and reward machines.

## 387 References

- 388 David Abel, Will Dabney, Anna Harutyunyan, Mark K. Ho, Michael L. Littman, Doina Precup, and  
389 Satinder Singh. On the expressivity of markov reward. In *NeurIPS*, pp. 7799–7812, 2021.
- 390 Saleema Amershi, Maya Cakmak, William Bradley Knox, and Todd Kulesza. Power to the people:  
391 The role of humans in interactive machine learning. *AI magazine*, 35(4):105–120, 2014.
- 392 Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. Con-  
393 crete problems in ai safety. *arXiv preprint arXiv:1606.06565*, 2016.
- 394 Andrea Bajcsy, Dylan P Losey, Marcia K O'Malley, and Anca D Dragan. Learning from physical  
395 human corrections, one feature at a time. In *Proceedings of the 2018 ACM/IEEE International*  
396 *Conference on Human-Robot Interaction*, pp. 141–149, 2018.

- 397 Kathy Baxter, Catherine Courage, and Kelly Caine. *Understanding your users: a practical guide to*  
398 *user research methods*. Morgan Kaufmann, 2015.
- 399 Erdem Biyik and Dorsa Sadigh. Batch active preference-based learning of reward functions. In  
400 *Conference on robot learning*, pp. 519–528. PMLR, 2018.
- 401 Serena Booth, W. Bradley Knox, Julie Shah, Scott Niekum, Peter Stone, and Alessandro Allievi.  
402 The perils of trial-and-error reward design: Misdesign through overfitting and invalid task speci-  
403 fications. In *AAAI*, pp. 5920–5929. AAAI Press, 2023.
- 404 Anthony Brohan, Yevgen Chebotar, Chelsea Finn, Karol Hausman, Alexander Herzog, Daniel Ho,  
405 Julian Ibarz, Alex Irpan, Eric Jang, Ryan Julian, et al. Do as i can, not as i say: Grounding  
406 language in robotic affordances. In *Conference on robot learning*, pp. 287–318. PMLR, 2023.
- 407 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
408 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
409 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 410 Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep  
411 reinforcement learning from human preferences. *Advances in neural information processing sys-*  
412 *tems*, 30, 2017.
- 413 Michael T Cox. A model of planning, action, and interpretation with goal reasoning. In *Proceedings*  
414 *of the 4th Annual Conference on Advances in Cognitive Systems*, pp. 48–63, 2016.
- 415 Nadia Figueroa, Salman Faraji, Mikhail Koptev, and Aude Billard. A dynamical system approach  
416 for adaptive grasping, navigation and co-manipulation with humanoid robots. In *2020 IEEE*  
417 *International conference on robotics and automation (ICRA)*, pp. 7676–7682. IEEE, 2020.
- 418 Hector Geffner and Blai Bonet. *A concise introduction to models and methods for automated plan-*  
419 *ning*. Morgan & Claypool Publishers, 2013.
- 420 Sandra G Hart. Nasa task load index (tlx). 1986.
- 421 Rodrigo Toro Icarte, Toryn Q Klassen, Richard Valenzano, and Sheila A McIlraith. Reward ma-  
422 chines: Exploiting reward function structure in reinforcement learning. *Journal of Artificial In-*  
423 *telligence Research*, 73:173–208, 2022.
- 424 W Bradley Knox and Peter Stone. Interactively shaping agents via human reinforcement: The tamer  
425 framework. In *Proceedings of the fifth international conference on Knowledge capture*, pp. 9–16,  
426 2009.
- 427 W Bradley Knox, Stephane Hatgis-Kessell, Serena Booth, Scott Niekum, Peter Stone, and Alessan-  
428 dro Allievi. Models of human preference for learning reward functions. *Transactions on Machine*  
429 *Learning Research*, 2022.
- 430 W. Bradley Knox, Alessandro Allievi, Holger Banzhaf, Felix Schmitt, and Peter Stone. Reward  
431 (Mis)design for autonomous driving. *Artificial Intelligence*, 316(103829), 2023.
- 432 Dylan P Losey and Marcia K O’Malley. Including uncertainty when learning from human correc-  
433 tions. In *Conference on Robot Learning*, pp. 123–132. PMLR, 2018.
- 434 James MacGlashan, Mark K Ho, Robert Loftin, Bei Peng, Guan Wang, David L Roberts, Matthew E  
435 Taylor, and Michael L Littman. Interactive learning from policy-dependent human feedback. In  
436 *International conference on machine learning*, pp. 2285–2294. PMLR, 2017.
- 437 Malek Mecherghi and Sarath Sreedharan. Goal alignment: Re-analyzing value alignment prob-  
438 lems using human-aware AI. In *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI*  
439 *2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024,*  
440 *Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2014, February*  
441 *20-27, 2024, Vancouver, Canada*, pp. 10110–10118. AAAI Press, 2024.

- 442 Martin L Puterman. Markov decision processes. *Handbooks in operations research and management*  
443 *science*, 2:331–434, 1990.
- 444 Harish Ravichandar, Athanasios S Polydoros, Sonia Chernova, and Aude Billard. Recent advances  
445 in robot learning from demonstration. *Annual review of control, robotics, and autonomous sys-*  
446 *tems*, 3(1):297–330, 2020.
- 447 Stuart Russell. Human-compatible artificial intelligence., 2022.
- 448 Rohin Shah, Vikrant Varma, Ramana Kumar, Mary Phuong, Victoria Krakovna, Jonathan Uesato,  
449 and Zac Kenton. Goal misgeneralization: Why correct specifications aren’t enough for correct  
450 goals. *arXiv preprint arXiv:2210.01790*, 2022.
- 451 Herbert A Simon. *The Sciences of the Artificial, reissue of the third edition with a new introduction*  
452 *by John Laird*. MIT press, 2019.
- 453 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- 454 Christopher Taylor. Aristotle on practical reason. In *The Oxford Handbook of Topics in Philosophy*.  
455 Oxford University Press, 2019. ISBN 9780199935314.
- 456 Andrea L Thomaz and Cynthia Breazeal. Teachable robots: Understanding human teaching behavior  
457 to build more effective robot learners. *Artificial Intelligence*, 172(6-7):716–737, 2008.
- 458 Nhi Tran. Goals vs. actions as user-facing representations for robot programming. In *2024 AAAI*  
459 *Fall Symposium Series*. AAAI, 2024.
- 460 Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul  
461 Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv*  
462 *preprint arXiv:1909.08593*, 2019.

## Supplementary Materials

*The following content was not necessarily subject to peer review.*

### 9 Raw Nasa TLX Scores

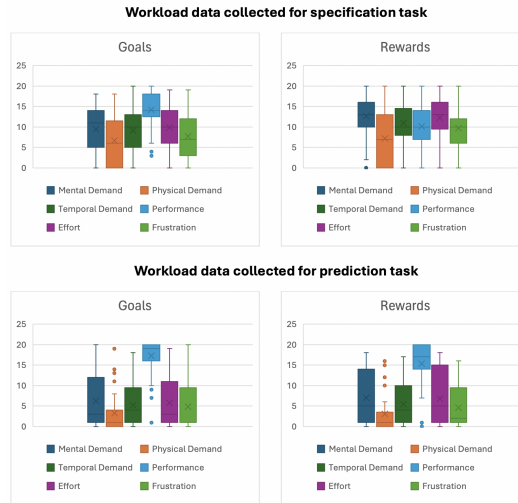


Figure 4: Box tables representing the NASA TLX score collected for both specification and prediction tasks.

### 10 Subjective Feedback from Participants

Here is the subjective feedback provided by the participants for each of the objective specification mechanisms. Here the participants were asked to select the objective mechanisms they felt most closely matched the description provided

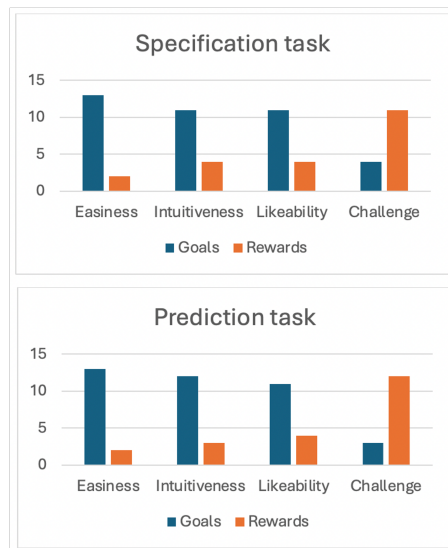


Figure 5: The raw number of selections provided by the participants for each task.



471 11 Screenshots from the Variants

Here are some of the screenshots from the two variants of the specification tasks.

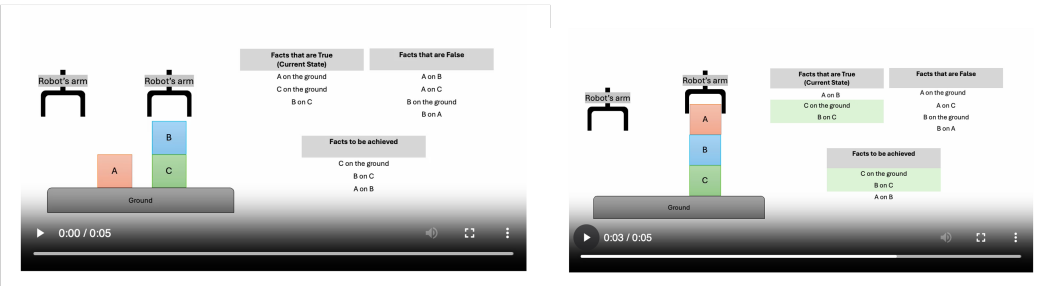


Figure 6: Screenshots from the new video for the first variant that highlights the false facts and how they change over actions.

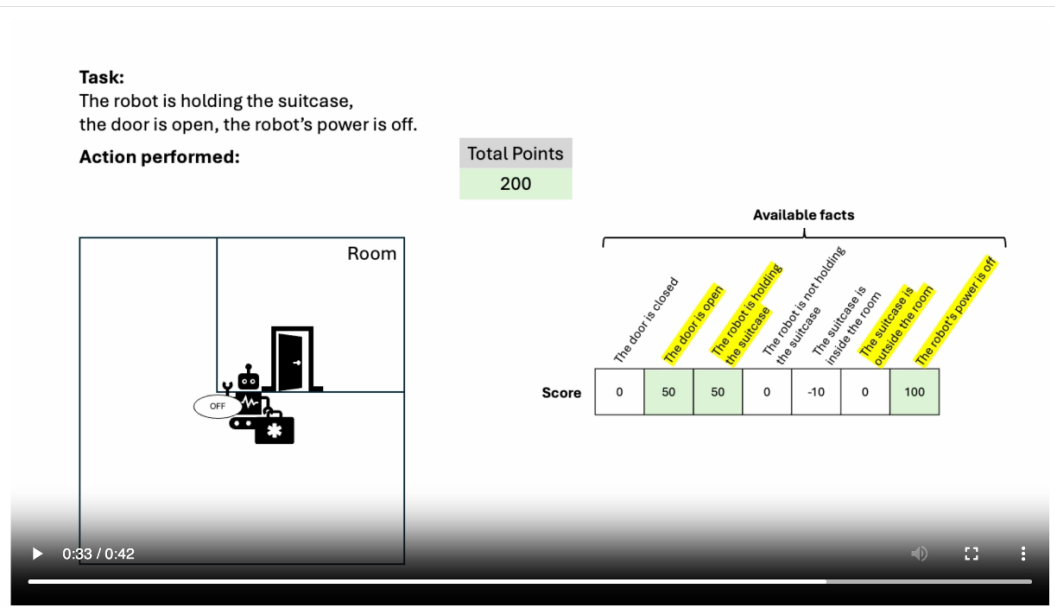


Figure 7: Screenshots from the second variant that shows the task and the new reward specification mechanism.