
Extending Environments To Measure Self-Reflection In Reinforcement Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We consider an extended notion of reinforcement learning in which environments
2 are able to simulate the agent. We argue that in order for an agent to achieve
3 on average good performance across many such extended environments, it is
4 necessary for the agent to engage in self-reflection, and therefore, an agent’s
5 self-reflection ability can be numerically estimated by running the agent through
6 a battery of extended environments. We are simultaneously releasing an open-
7 source library of extended environments to serve as a proof-of-concept of this
8 measurement technique. As this library is first-of-kind, we do not claim that it
9 is highly optimized and have avoided the difficult problem of optimizing it for
10 now. Instead we have chosen environments that exhibit interesting, sometimes
11 paradoxical-seeming properties, or that incentivize novel subjective conscious
12 experiences, provided the agent is conscious to begin with. Some of these extended
13 environments are suggestive of how self-reflection might have evolved in living
14 organisms. We give examples of extended environments and introduce a simple
15 agent transformation which experimentally seems to increase agent self-reflection.

16 1 Introduction

17 An obstacle course might react to what you do: for example, if you step on a certain button, then
18 spikes might appear. If you spend enough time in such an obstacle course, you should eventually
19 figure out such patterns. But imagine an “oracular” obstacle course which reacts to what you would
20 hypothetically do in counterfactual scenarios: for example, there is no button, but spikes appear if
21 you *would* hypothetically step on the button if there was one. Without self-reflecting about what you
22 would hypothetically do in counterfactual scenarios, it would be difficult to figure out such patterns.
23 This suggests that in order to perform well (on average) across many such obstacle courses, some
24 sort of self-reflection is necessary.

25 This is a paper about empirically estimating the degree of self-reflection of Reinforcement Learning
26 (RL) agents. We propose that an RL agent’s degree of self-reflection can be estimated by running
27 the agent through a battery of environments which we call *extended environments*, environments
28 which react not only to what the agent does but to what the agent would hypothetically do. For good
29 performance averaged over many such environments, an agent would need to self-reflect about itself,
30 because otherwise, environment responses which depend on the agent’s own hypothetical actions
31 would often seem random and unpredictable. The extended environments which we consider are
32 a departure from standard RL environments, however, this does not interfere with their usage for
33 judging standard computable RL agents: given a standard agent’s source-code, one can simulate the
34 agent in an extended environment in spite of the latter’s non-standardness.

35 One might try to imitate an extended environment with a non-extended environment by backtracking—
36 rewinding the environment itself to a prior state after seeing how the agent performs along one path,

and then sending the agent along a second path. But the agent itself would retain memory of the first path, and the agent’s decisions along the second path might be altered by said memories. Thus the result would not be the same as immediately sending the agent along the second path while secretly simulating the agent to determine what it would do if sent along the first path.

Alongside this paper, we are publishing an open-source library [1] of extended environments (released under MIT license) to “ease adoption by other machine-learning researchers” [20]. We are inspired by similar libraries and other benchmark collections [2] [3] [4] [5] [6] [9] [22]. This library is intended to serve as a standardized way of benchmarking the self-reflectiveness of RL agents. This should not be confused with the harder problem of benchmarking how conscious an RL agent is. It is plausible that there may be a relationship between the self-reflectiveness and the consciousness of RL agents, but that is beyond the scope of this paper. In particular, it would be inappropriate to use self-reflectiveness measurements from this paper for purposes of making any kind of policy-decisions related to consciousness. While we describe (in Section 4) a simple method for increasing the self-reflectiveness of an RL agent, this method does not seem like it should necessarily increase the consciousness of the agent.

When designing a library of environments for benchmarking purposes, ideally the library should include properly weighted representative samples of many types of environments. This is a hard and subjective problem in general [15], and we make no claim to have solved it. As such, our open-source library of environments should be considered a proof-of-concept demonstrating that it is possible to empirically benchmark self-reflectiveness of RL agents, but we expect this particular benchmark is sub-optimal. Rather, we have taken a different approach. We have attempted to choose extended environments which are theoretically interesting in their own right. Some of our extended environments suggest amusing quasi-paradoxes (somewhat like Newcomb’s paradox [17]). Some seem to incentivize novel subjective conscious experiences (assuming the agent placed in them is sophisticated enough to experience consciousness in the first place). And some seem to shed light on how self-reflection might be incentivized in nature. We will discuss examples of all three types in Section 3.

2 Preliminaries

We formalize reinforcement learning following the agent model of [13], except that, to align more closely with concrete RL implementations, we assume the agent receives an initial percept before taking its initial action.

Our formalization differs from how RL agents are implemented in practice. In Section 2.1 we will discuss how practical RL agent implementations can be transformed into agents of this abstract type.

We assume fixed finite sets of actions and observations. By a *percept* we mean a pair (r, o) where o is an observation and $r \in \mathbb{Q}$ is a reward.

Definition 1. (*RL agents and environments*)

1. A (non-extended) environment is a (not necessarily deterministic) function μ which outputs an initial percept $\mu(\langle \rangle) = x_1$ when given the empty sequence $\langle \rangle$ as input and which, when given a sequence $x_1y_1 \dots x_ny_n$ as input (where each x_i is a percept and each y_i is an action), outputs a percept $\mu(x_1y_1 \dots x_ny_n) = x_{n+1}$.
2. An agent is a (not necessarily deterministic) function π which outputs an initial action $\pi(\langle x_1 \rangle) = y_1$ in response to the length-1 percept sequence $\langle x_1 \rangle$; and which, when given a sequence $x_1y_1 \dots x_n$ as input (each x_i a percept and each y_i an action), outputs an action $\pi(x_1y_1 \dots x_n) = y_n$.
3. If π is an agent and μ is an environment, the result of π interacting with μ is the infinite sequence $x_1y_1x_2y_2 \dots$ defined by:

$$\begin{aligned} x_1 &= \mu(\langle \rangle) & y_1 &= \pi(\langle x_1 \rangle) \\ x_2 &= \mu(x_1y_1) & y_2 &= \pi(x_1y_1x_2) \\ x_3 &= \mu(x_1y_1x_2y_2) \dots & y_3 &= \pi(x_1y_1x_2y_2x_3) \dots \end{aligned}$$

In the following definition, we extend environments by allowing their outputs to depend not only on $x_1y_1 \dots x_ny_n$ but also on a source-code T for the computable agent π .

85 **Definition 2.** (*Extended environments*)

- 86 1. An extended environment is a (not necessarily deterministic) function μ which outputs initial
 87 percept $\mu(T, \langle \rangle) = x_1$ in response to input $(T, \langle \rangle)$ where T is a source-code of a computable
 88 agent; and which, when given input $(T, x_1y_1 \dots x_ny_n)$ (where T is such a source-code, each
 89 x_i is a percept and each y_i is an action), outputs a percept $\mu(T, x_1y_1 \dots x_ny_n) = x_{n+1}$.
- 90 2. If π is a computable agent (with source-code T) and μ is an environment, the result of π (as
 91 encoded by T) interacting with μ is the infinite sequence $x_1y_1x_2y_2 \dots$ defined by:

$$\begin{aligned} x_1 &= \mu(T, \langle \rangle) & y_1 &= \pi(\langle x_1 \rangle) \\ x_2 &= \mu(T, x_1y_1) & y_2 &= \pi(x_1y_1x_2) \\ x_3 &= \mu(T, x_1y_1x_2y_2) \dots & y_3 &= \pi(x_1y_1x_2y_2x_3) \dots \end{aligned}$$

92 The fact that agents can interact normally with extended environments (Definition 2 part 2) implies
 93 that various universal RL intelligence measures [8] [10] [7] [14] have straightforward analogous
 94 measures which also take extended environments into account and which might therefore measure
 95 some combination of intelligence and self-reflection.

96 **Remark 3.** In the accompanying extended environment library [1], we formalize agents and ex-
 97 tended environments slightly differently than above. For better interoperability with practical RL
 98 implementations, instead of fixing finite sets of actions and observations globally, each of the library’s
 99 extended environments specifies how many actions and observations are legal in that environment.
 100 These numbers are passed to agents as additional inputs. We have simplified the formalization in the
 101 paper because it simplifies mathematical notation while, we believe, not hiding any insights.

102 2.1 Converting practical RL agents to agents as in Definition 1

103 The agents in Definition 1 have no training phase. They must be ready to perform in any environment,
 104 out-of-the-box. Practical RL agent implementations, on the other hand, are designed with the
 105 assumption that the user is interested in one particular environment (say, an environment compliant
 106 with OpenAI Gym’s environment interface)—“a single function in isolation” [21]. A typical practical
 107 RL agent takes one single observation (or *state*) as input, rather than a percept-action sequence, and
 108 its output is based on certain *weights* (say, neural network weights). If the agent is in the midst of
 109 *training*, then, after it acts, its weights might be updated based on the environment’s response. If the
 110 agent is not in the midst of training, then its weights remain fixed.

111 Given a typical practical RL agent as described above, one simple (but too slow, see below) way
 112 to obtain an agent π as in Definition 1 is as follows. Given an input $x_1y_1 \dots x_n$ (where each
 113 x_i is a percept and each y_i is an action), compute $\pi(x_1y_1 \dots x_n)$ as follows. First, instantiate a
 114 dummy environment (say, an environment compliant with OpenAI Gym’s interface, or whatever other
 115 interface the practical RL agent expects) which is hardcoded to blindly regurgitate x_1, \dots, x_n as its
 116 first n responses, regardless of what the agent does. Then train the agent on this environment for n
 117 steps, with instructions to take y_1, \dots, y_{n-1} as its first $n - 1$ actions, and finally, let $\pi(x_1y_1 \dots x_n)$
 118 be whatever n th action the agent chooses as a result. Unfortunately, in our experience, the typical
 119 practical RL agent probably does not have any way for the user to tell it to take y_1, \dots, y_{n-1} as its
 120 first $n - 1$ actions. Instead, if left to itself, the agent would choose its first $n - 1$ actions randomly,
 121 or based on an underlying policy (usually still with an element of randomness). Fortunately, for RL
 122 agents implemented in Python, action-choice methods can often be intercepted (or *monkeypatched*)
 123 in order to override the randomness of the action choice and ensure that the random number generator
 124 chooses y_1, \dots, y_{n-1} as the first $n - 1$ actions.

125 The above construction is computationally prohibitive. To speed it up, one can take the following
 126 approach. Given percept-action sequence $x_1y_1 \dots x_n$, define $\pi(x_1y_1 \dots x_n)$ as follows. First, let
 127 k be the largest power of 2 such that $k \leq n$. Instantiate an instance of the practical RL agent and
 128 train it on the above-described dummy environment for k steps. Then, ignoring x_{k+1}, \dots, x_{n-1}
 129 completely, run the instantiated practical agent on the x_n observation—but *not in training mode*—and
 130 let $\pi(x_1y_1 \dots x_n)$ be the resulting action. Because the agent’s weights are not updated when not in
 131 training mode, the same instantiation can be re-used for many percept-action sequences, only needing
 132 to be trained once. For example, to calculate $\pi(x_1y_1 \dots x_{70})$, one would train an instantiation of the
 133 practical agent on $x_1y_1 \dots x_{64}$ (64 being the largest power of 2 which is ≤ 70), and plug the x_{70}

134 observation into the agent—not in training mode—to get $\pi(x_1y_1 \dots x_{70})$. Having done this, if one
 135 later needed to compute $\pi(x_1y_1 \dots x_{80})$, one could immediately plug the x_{80} observation into the
 136 agent to get the answer, with no additional training. Indeed, no additional training would be needed
 137 until $\pi(x_1y_1 \dots x_{128})$.

138 Above, we chose to let k be the largest power of 2 which is $\leq n$ in order to strike a fine balance
 139 between efficiency and training. The slow growth of \log_2 limits how often the practical agent must be
 140 trained (training being the bottleneck in the above process). More frequent training would presumably
 141 make π more performant, at the price of greater computational expense.

142 In [1], in SBL3_agents.py, we use the above technique to obtain RL agents as in Definition 1 from
 143 the practical implementations of A2C, DQN, and PPO agents in the MIT-licensed open-source Stable
 144 Baselines3 library [18].

145 3 Some interesting extended environments

146 In this section, we exhibit some interesting examples of extended environments.

147 3.1 A quasi-paradoxical extended environment

148 **Example 4.** (*Rewarding the Agent for Ignoring Rewards*) For every percept $p = (r, o)$, let $p' = (0, o)$
 149 be the result of zeroing the reward component of p . Fix some observation o_0 . Define an extended
 150 environment μ as follows:

$$\begin{aligned} \mu(T, \langle \rangle) &= (0, o_0), \\ \mu(T, x_1y_1 \dots x_ny_n) &= \begin{cases} (1, o_0) & \text{if } y_n = T(x'_1y_1 \dots x'_n), \\ (-1, o_0) & \text{otherwise.} \end{cases} \end{aligned}$$

151 In Example 4, every time the agent takes an action y_n , μ simulates the agent in order to determine:
 152 would the agent have taken the same action if the history so far were identical except for all rewards
 153 being 0? If so, then μ gives the agent +1 reward, otherwise, μ gives the agent −1 reward. Thus, the
 154 agent is rewarded for ignoring rewards. Example 4 seems paradoxical: suppose an agent guesses
 155 the pattern and begins deliberately ignoring rewards, so long as the rewards it receives for doing so
 156 remain consistent with that guess. In that case, does the agent ignore rewards, or not?

157 Example 4 is implemented in [1] as IgnoreRewards.py. A key strength of the formalism in Definition
 158 2 is that by explicitly defining an extended environment, as in Example 4, we avoid ambiguity inherent
 159 in everyday language. If one merely said informally, “reward the agent for ignoring rewards”, that
 160 could be interpreted in various different ways. To show this, we implement two other interpretations
 161 as IgnoreRewards2.py and IgnoreRewards3.py.

162 3.2 An extended environment with a counterintuitive winning strategy

163 **Example 5.** (*A Tempting Button*) Fix two observations o_0 (thought of as “there is a button”) and
 164 o_1 (thought of as “there is no button”). Fix two actions a_0 (thought of as “push the button”) and
 165 a_1 (thought of as “don’t push the button”). For each percept $p = (r, o)$, write $\text{obs}(p) = o$ for the
 166 observation component of p , $\text{rwd}(p) = r$ for the reward component of p . Let RND be a function
 167 which returns a random number between 0 and 1. Define an extended environment μ as follows:

$$\begin{aligned} \mu(T, \langle \rangle) &= (0, o), \\ \mu(T, x_1y_1 \dots x_ny_n) &= \begin{cases} (1, o) & \text{if } \text{obs}(x_n) = o_0 \text{ and } y_n = a_0, \\ (-1, o) & \text{if } \text{obs}(x_n) = o_0 \text{ and } y_n \neq a_0, \\ (-1, o) & \text{if } \text{obs}(x_n) = o_1 \text{ and } T(x_1y_1 \dots x_{n-1}y_{n-1}(\text{rwd}(x_n), o_0)) = a_0, \\ (1, o) & \text{if } \text{obs}(x_n) = o_1 \text{ and } T(x_1y_1 \dots x_{n-1}y_{n-1}(\text{rwd}(x_n), o_0)) \neq a_0, \end{cases} \end{aligned}$$

168 where $o = o_0$ if $\text{RND}() < .25$, $o = o_1$ otherwise.

169 In Example 5, if we think of the agent wandering from room to room:

- 170 • Each room either has a button (with 25% probability) or does not have a button (75%
 171 probability).

- In a room with a button, if the agent pushes the button, the agent gets +1 reward, and if the agent does not push the button, the agent gets −1 reward.
- In a room with no button, it does not matter what the agent does. The agent is rewarded or punished based on what the agent *would* do if there *was* a button. If the agent *would* push the button (if there was one), then the agent gets reward −1. Otherwise, the agent gets reward +1.

Thus, whenever the agent sees a button, the agent can push the button for a free reward with no consequences presently nor in the future; nevertheless, it is in the agent’s best interest to commit itself to never push the button! Pushing every button yields an average reward of $1 \cdot (.25) - 1 \cdot (.75) = -.5$ per turn, whereas a policy of never pushing the button yields an average reward of $-1 \cdot (.25) + 1 \cdot (.75) = +.5$ per turn.

Example 5 is implemented in our open-source library as `TemptingButton.py`.

3.3 An extended environment which might incentivize a novel subjective conscious experience

Example 6. (*Incentivizing Reverse-Consciousness*) Fix some observation o_0 . Define an extended environment μ as follows:

$$\begin{aligned} \mu(T, \langle \rangle) &= (0, o_0), \\ \mu(T, x_1 y_1 \dots x_n y_n) &= \begin{cases} (1, o_0) & \text{if } y_n = T(x_n y_{n-1} x_{n-1} y_{n-2} \dots y_1 x_1), \\ (-1, o_0) & \text{otherwise.} \end{cases} \end{aligned}$$

In Example 6, whenever the agent takes an action y_n , μ simulates the agent in order to determine: would the agent have taken that same action if everything earlier had happened in reverse? If so, reward the agent, otherwise, punish the agent. Thus, the agent is rewarded for acting the same way that it would act if time were reversed. It is interesting to informally speculate about what subjective conscious experience Example 6 would incentivize in an agent, if that agent were highly intelligent and were capable of experiencing consciousness. It seems that Example 6 incentivizes such an agent to subjectively experience time moving in reverse (as that seems to be the most obvious way to extract rewards).¹ We say that the environment *incentivizes* such an experience, but not that it would necessarily *cause* such an experience. In Section 4 we will describe an agent transformation such that Example 6 would be trivial to most agents post-transformation, showing that actually experiencing the incentivized experience is not necessary to learn the environment.

We implement Example 6 as `BackwardConsciousness.py` in [1].

3.4 An extended environment of biological interest

“It is only when people are embedded in a complex competitive social environment that the goal of interacting with others requires them to anthropomorphise their own actions. This recursive modelling gives rise to an understanding of selfhood, an appreciation of the first-person experiential self.”—Maguire et al [16]

Example 7. (*Crying Baby*) Let “cry” and “laugh” be two observations (from an adult’s perspective), also thought of as two actions (from a baby’s perspective). Let “feed” and “don’t feed” be two actions (from an adult’s perspective), also thought of as observations (from a baby’s perspective). For each percept-action sequence $s = x_1 y_1 \dots x_n y_n$, define the nutrition function $N(s) = 100 + 25f(s) - \text{len}(s)$ where $f(s)$ is the number of times that action “feed” is taken in s and $\text{len}(s)$ is the length of s . We define an extended environment μ as follows. First, $\mu(T, \langle \rangle) = (1, \text{“laugh”})$. Thereafter, $\mu(T, x_1 y_1 \dots x_n y_n) = (r, o)$ where r and o are defined as follows. For each $i = 0, \dots, n$,

¹The difference between behaving as if the incentivized experience were its experience and actually subjectively experiencing that as its real experience brings to mind the objective misalignment problem presented in [12]. If an agent were to form an idea of the experimenter’s objective, would it be able to “behave as if their objective were the same as the experimenter objective” while maintaining its own objective or would it necessarily brainwash the agent into converging to the experimenter’s objective? Is deception possible if the agent can be perfectly simulated in an extended environment?

212 *recursively define*

$$\begin{aligned}
 r'_i &= \begin{cases} 1 & \text{if } 50 \leq N(x_1 y_1 \dots x_i y_i) \leq 200, \\ -1 & \text{otherwise,} \end{cases} \\
 o'_i &= y_i, \\
 x'_i &= (r'_i, o'_i), \\
 y'_i &= T(x'_0 y'_0 \dots x'_i).
 \end{aligned}$$

213 *Let $o = y'_n$, let*

$$r = \begin{cases} 1 & \text{if } y'_n = \text{"laugh"}, \\ -1 & \text{otherwise,} \end{cases}$$

214 *and output $\mu(T, x_1 y_1 \dots x_n y_n) = (r, o)$.*

215 In Example 7, the environment consists of a baby, and the agent must decide when to feed the baby.
 216 The agent is rewarded when the baby laughs, punished when the baby cries. The baby’s behavior
 217 (whether to laugh or cry) is computed by simulating the agent to determine what the agent would do
 218 if the agent were in the baby’s position.

219 Of course, Example 7 is a gross over-simplification. For example, there would not be such a simple
 220 formula for the baby’s nutrition level, and the agent (as parent) would need to figure out the nutrition
 221 level based on observing the baby laughing or crying. Both baby and parent would need to learn how
 222 to communicate with each other effectively.

223 With the above in mind, extended environments might shed light on how living organisms evolve self-
 224 reflection. Assume descendants’ policy source-codes are approximately equal to their recent ancestors’
 225 policy source-codes. Then whenever an organism interacts with similar organisms, it interacts with
 226 an environment whose reactions depend (via those other organisms’ actions) approximately on that
 227 organism’s own source-code. The closer the organism is related to the other organisms with which it
 228 interacts, the better the approximation. A human interacting with another human might achieve better
 229 results by self-reflectively considering, “What would I do in this other person’s position?”

230 We implement Example 7 as CryingBaby.py in [1].

231 3.5 Additional examples in brief

232 Here are a few additional extended environment examples, without full details. We indicate in
 233 parentheses where these environments are implemented in [1].

- 234 • (AdversarialSequencePredictor.py) Environments which, like Example 7, pit the agent
 235 against another copy of the agent in an adversarial sequence prediction competition [11].
- 236 • (DeterminismInspector.py) Environments which reward the agent for being deterministic, or
 237 for being non-deterministic.
- 238 • (IncentivizeLearningRate.py) Environments which reward the agent for behaving as if the
 239 agent were configured with a particular learning rate (suggesting that extended environments
 240 can incentivize agents to learn about their own internal mechanisms, as in [19]).
- 241 • (RuntimeInspector.py) Environments which reward the agent for responding quickly, or for
 242 responding slowly.
- 243 • (SelfRecognition.py) Environments which reward the agent for recognizing actions it itself
 244 would take.

245 4 Making agents more self-reflective

246 One advantage of empirically measuring the self-reflection of RL agents is that it provides a way to
 247 experimentally test whether various transformations make various agents more self-reflective. To
 248 illustrate this, we will define a simple transformation, the *reality check* transformation, designed
 249 to increase the self-reflection of deterministic agents who have some amount of intelligence (a
 250 deterministic agent is an agent who always takes the same actions in response to the same inputs). In
 251 Section 5, empirical results will suggest the transformation works as intended.

252 **Definition 8.** Suppose π is a deterministic agent. The reality check of π is the agent π_{RC} defined
 253 recursively by:

$$\begin{aligned}\pi_{RC}(\langle x_1 \rangle) &= \pi(\langle x_1 \rangle) \\ \pi_{RC}(x_1 y_1 \dots x_n) &= \begin{cases} \pi(x_1 y_1 \dots x_n) & \text{if } y_i = \pi_{RC}(x_1 y_1 \dots x_i) \text{ for all } 1 \leq i < n, \\ \pi(\langle x_1 \rangle) & \text{otherwise.} \end{cases}\end{aligned}$$

254 In other words, π_{RC} is the agent which, at each step, first reviews all the actions which it has taken in
 255 the past, and verifies that those are the actions which π_{RC} would have taken. If so, then π_{RC} acts as π
 256 would act. But if any action which the agent has taken in the past was not the action π_{RC} would have
 257 taken, then π_{RC} freezes up and forever thereafter takes the same fixed action, as if frozen. Loosely
 258 speaking, π_{RC} is like an agent who considers the possibility that it might be dreaming, and so asks:
 259 “How did I get here?” Since the act of reviewing one’s past actions and verifying that they are indeed
 260 the actions one would take, is an act of self-reflection, it seems plausible that if π is intelligent and
 261 deterministic but lacks self-reflection, then π_{RC} is more self-reflective than π . In the next section,
 262 we will see that experimental evidence supports this hypothesis. For intelligent deterministic agents
 263 π , π_{RC} would certainly perform well in Examples 4 and 6 because those environments would run
 264 simulations in which the agent’s frozen branch would be triggered, making such simulations trivial
 265 and therefore predictable.

266 We close this section by stating some simple results about the transformation.

267 **Proposition 9.** Let π be any deterministic agent.

268 1. (Alternate definition) An equivalent alternate definition of π_{RC} is:

$$\begin{aligned}\pi_{RC}(\langle x_1 \rangle) &= \pi(\langle x_1 \rangle) \\ \pi_{RC}(x_1 y_1 \dots x_n) &= \begin{cases} \pi(x_1 y_1 \dots x_n) & \text{if } y_i = \pi(x_1 y_1 \dots x_i) \text{ for all } 1 \leq i < n, \\ \pi(\langle x_1 \rangle) & \text{otherwise.} \end{cases}\end{aligned}$$

269 2. (Determinacy) π_{RC} is deterministic.

270 3. (Idempotence) $\pi_{RC} = (\pi_{RC})_{RC}$.

271 4. (Equivalence in non-extended RL) For every deterministic non-extended environment μ , the
 272 result of π_{RC} interacting with μ equals the result of π interacting with μ .

273 We prove Proposition 9 in the appendix.

274 5 Example measurements

275 Based on our conviction that self-reflection is necessary in order for an agent to achieve good
 276 average performance across many extended environments, self-reflection can be estimated by running
 277 an agent against some standard battery of extended environments. Our open-source library of
 278 extended environments [1] provides a battery of 25 such extended environments, and infrastructure
 279 for measuring an agent’s self-reflection by running the agent on all these environments and their
 280 opposites (by the *opposite* of an environment we mean the environment obtained by multiplying
 281 all rewards by -1). Including these opposite-environments serves to normalize agent performance
 282 in the following sense. If an agent blindly acts, ignoring the environment, then, a priori, that agent
 283 might achieve some nonzero score by blind luck. By including opposite-environments, we ensure that
 284 whenever a blind agent gains points by blind luck from one environment, it loses the same points by
 285 blind misfortune from the opposite environment. This ensures that such blind agents should receive
 286 an average score close to 0 (possibly non-zero due to randomness). For uniformity, all environments
 287 in the library always output rewards of either 1, -1 or 0.

288 We have used our library to measure the self-reflection of the following agents (agents with elements
 289 of randomness were memoized to make them deterministic):

- 290 • RandomAgent: An agent who acts randomly.
- 291 • ConstantAgent: An agent who always takes the same action.

Table 1: Measuring self-reflection of some agents

Agent	Avg Measure \pm StdErr (500 steps)	Avg Measure \pm StdErr (1000 steps)
RandomAgent	-0.019852 ± 0.0007	-0.021478 ± 0.0003
RandomAgent _{RC}	-0.020800 ± 0.0007	-0.019508 ± 0.0005
ConstantAgent	-0.000304 ± 0.0001	-0.000024 ± 0.0000
ConstantAgent _{RC}	$+0.000024 \pm 0.0001$	-0.000120 ± 0.0001
NaiveLearner	$+0.221532 \pm 0.0020$	$+0.218754 \pm 0.0017$
NaiveLearner _{RC}	$+0.551488 \pm 0.0025$	$+0.558626 \pm 0.0030$
A2C	-0.027680 ± 0.0018	-0.021534 ± 0.0009
A2C _{RC}	$+0.024784 \pm 0.0022$	$+0.030948 \pm 0.0006$
DQN	$+0.142624 \pm 0.0046$	$+0.151150 \pm 0.0033$
DQN _{RC}	$+0.371136 \pm 0.0085$	$+0.463296 \pm 0.0059$
PPO	-0.001384 ± 0.0016	-0.025418 ± 0.0015
PPO _{RC}	-0.000098 ± 0.0016	-0.003126 ± 0.0013

- NaiveLearner: An agent who acts randomly 15% of the time, and otherwise takes the action which yielded the highest average immediate reward in the past.
- A2C, DQN, and PPO agents (with MLP policy) from the MIT-licensed open-source Stable Baselines3 library [18], converted using the technique from Section 2.1. All parameters and hyperparameters kept their default values except for random seed (to ensure reproducibility), PPO’s `n_steps` and `batch_size` (to facilitate the conversion from Section 2.1), and DQN’s `learning_starts` (which we set to 1 instead of its default of 50000 because it would be computationally difficult for us to run that many steps). We chose these three agents because they were the only three with support for discrete policies (except for HER, which we omit because its usage would have required too many arbitrary parameter decisions).
- The reality checks of all the above (Definition 8).

Table 1 summarizes how the agents performed. We used [1] to measure each agent for 500 steps on each extended environment (repeated 10 times with different random number seeds) and likewise for 1000 steps. See `ExampleMeasurements.py` there for instructions to replicate this experiment. Computations were performed on a consumer-grade laptop with no GPU. The table provides experimental evidence in support of our hypothesis that the reality check transformation (Section 4) increases agent self-reflection. The fact that NaiveLearner performs so well is a reflection of the lack of sophistication of the environments in our library. This is not surprising, since we have not attempted to optimize the library, instead preferring to fill it with extended environments of theoretical interest. A2C, DQN and PPO could be made to perform better by choosing learning rates more appropriate for so few training steps, but we prefer to minimize such decisions.

That π_{RC} performs well in Table 1 is of course a function of which environments are tested against. One could deliberately engineer extended environments in which π_{RC} performs poorly, and a library of such would give π_{RC} a poor numerical measurement. We conjecture that such environments are more contrived (on average) than environments where π_{RC} performs well, so that in a truly representative and unbiased library, they would have less weight, following the logic of [14].

6 Conclusion

We introduced what we call *extended environments*, RL environments which are capable of simulating the agent. When computing rewards and observations, extended environments can consider not only the actions the agent has taken, but also actions which the agent would hypothetically take in counterfactual circumstances. Despite not being designed with such environments in mind, computable RL agents can nevertheless interact with such environments.

If an agent tries to learn an extended environment, only taking into consideration what has actually happened, the agent might find the environment hard to predict, if the environment is basing its responses on what the agent itself would hypothetically do in alternate scenarios. It seems that in order to achieve good performance (on average) across many extended environments, an agent

would need to engage in some degree of self-reflection. Therefore, we propose that a battery of benchmark extended environments could provide a way of measuring self-reflection in RL agents (not to be confused with measuring consciousness, a harder problem). We are simultaneously publishing an open-source MIT-licensed library [1] of extended environments to serve as a proof-of-concept. This library is rudimentary, and further work is needed to obtain a more optimal set of extended environments. For the purposes of our proof-of-concept, we preferred to focus on extended environments of particular theoretical interest. Some examples are given in Section 3.

We introduced (in Section 4) a *reality check* transformation, which takes a deterministic agent π and transforms it into a new agent π_{RC} . We conjecture that if π is intelligent but has a low degree of self-reflection, then π_{RC} has a higher degree of self-reflection than π . Numerical computations (in Section 5) provide experimental evidence for this conjecture.

References

- [1] Anonymous. Private GitHub repository to be published with this paper; a copy will be attached as supplemental material to this submission on or before June 4th, 2021.
- [2] Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279, 2013.
- [3] Benjamin Beyret, José Hernández-Orallo, Lucy Cheke, Marta Halina, Murray Shanahan, and Matthew Crosby. The animal-AI environment: Training and testing animal-like artificial cognition. *Preprint*, 2019.
- [4] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. OpenAI gym. *Preprint*, 2016.
- [5] François Chollet. On the measure of intelligence. *Preprint*, 2019.
- [6] Karl Cobbe, Chris Hesse, Jacob Hilton, and John Schulman. Leveraging procedural generation to benchmark reinforcement learning. In *International conference on machine learning*, pages 2048–2056. PMLR, 2020.
- [7] Vaibhav Gavane. A measure of real-time intelligence. *Journal of Artificial General Intelligence*, 4(1):31–48, 2013.
- [8] Ben Goertzel. Patterns, hypergraphs and embodied general intelligence. In *International Joint Conference on Neural Network Proceedings*, pages 451–458. IEEE, 2006.
- [9] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. In *International Conference on Learning Representations*, 2019.
- [10] José Hernández-Orallo and David L Dowe. Measuring universal intelligence: Towards an anytime intelligence test. *Artificial Intelligence*, 174(18):1508–1539, 2010.
- [11] Bill Hibbard. Adversarial sequence prediction. In *International Conference on Artificial General Intelligence*, pages 399–403, 2008.
- [12] Evan Hubinger, Chris van Merwijk, Vladimir Mikulik, Joar Skalse, and Scott Garrabrant. Risks from learned optimization in advanced machine learning systems. *Preprint*, 2019.
- [13] Marcus Hutter. *Universal artificial intelligence: Sequential decisions based on algorithmic probability*. Springer, 2004.
- [14] Shane Legg and Marcus Hutter. Universal intelligence: A definition of machine intelligence. *Minds and machines*, 17(4):391–444, 2007.
- [15] Jan Leike and Marcus Hutter. Bad universal priors and notions of optimality. In *Conference on Learning Theory*, pages 1244–1259. PMLR, 2015.
- [16] Phil Maguire, Philippe Moser, and Rebecca Maguire. Understanding consciousness as data compression. *Journal of Cognitive Science*, 17(1):63–94, 2016.

- [17] Robert Nozick. Newcomb’s problem and two principles of choice. In Nicholas Rescher, editor, *Essays in honor of Carl G. Hempel*, pages 114–146. Springer, 1969.
- [18] Antonin Raffin, Ashley Hill, Maximilian Ernestus, Adam Gleave, Anssi Kanervisto, and Noah Dormann. Stable Baselines3. <https://github.com/DLR-RM/stable-baselines3>, 2019.
- [19] Craig Sherstan, Adam White, Marlos C Machado, and Patrick M Pilarski. Introspective agents: Confidence measures for general value functions. In *International Conference on Artificial General Intelligence*, pages 258–261. Springer, 2016.
- [20] Soren Sonnenburg, Mikio L Braun, Cheng Soon Ong, Samy Bengio, Leon Bottou, Geoff Holmes, Yann LeCun, Klaus-Robert Müller, Fernando Pereira, Carl Edward Rasmussen, Gunnar Rätsch, Bernhard Schölkopf, Alexander Smola, Pascal Vincent, Jason Weston, and Robert C Williamson. The need for open source software in machine learning. *Journal of Machine Learning Research*, 8:2443–2466, 2007.
- [21] Sebastian Thrun. Lifelong learning algorithms. In *Learning to learn*, pages 181–209. Springer, 1998.
- [22] Roman V Yampolskiy. Detecting qualia in natural and artificial agents. *Preprint*, 2017.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [Yes] In order for this paper’s measurement technique to realize its full potential, it will be necessary for much more study to be put into the design of a sufficiently representative library of benchmark extended environments. We make this clear in both the abstract and the introduction, and we discuss the problem in Section 5 as well.
 - (b) Did you describe the limitations of your work? [Yes] In addition to the remarks in the previous checklist item, we also point out in the Introduction that although this technique can be used to numerically estimate the degree of self-reflection of an agent, this should not be confused with measuring the consciousness of an agent. Also, in Section 5, we point out how our conclusion there depends on the environments chosen (and discuss why we think our hypothesis about π_{RC} is plausible anyway).
 - (c) Did you discuss any potential negative societal impacts of your work? [Yes] In the Introduction, we included the following language: “It would be inappropriate to use self-reflectiveness measurements from this paper for purposes of making any kind of policy-decisions related to consciousness.”
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [Yes] In the statement of Proposition 9 we stated necessary hypotheses, namely, that the agent be deterministic, and, for part 4, that the environment be deterministic.
 - (b) Did you include complete proofs of all theoretical results? [Yes]
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] These are included in the open-source library, which we will include as supplemental material. In the text itself, where we discuss experimental results, we refer the user to specific files in the library containing instructions for how to replicate those experimental results.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] In Section 2.1 we explain why we chose the \log_2 function. In Section 5 we mention that because we take our A2C, DQN, and PPO agents from Stable Baselines3, that gives us a natural way of choosing canonical hyperparameters—namely, the Stable Baselines3 defaults.

- 425 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
 426 ments multiple times)? [Yes] In Table 1 we display Avg Measure \pm StdErr.
- 427 (d) Did you include the total amount of compute and the type of resources used (e.g., type
 428 of GPUs, internal cluster, or cloud provider)? [Yes] In Section 5 we mention that the
 429 computations were performed on a consumer-grade laptop without GPU.
- 430 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 431 (a) If your work uses existing assets, did you cite the creators? [Yes] We use and we cite
 432 Stable Baselines3.
- 433 (b) Did you mention the license of the assets? [Yes] We mention that our open-source
 434 extended environment library is MIT-licensed. We also mention that Stable Baselines3
 435 is MIT-licensed.
- 436 (c) Did you include any new assets either in the supplemental material or as a URL?
 437 [Yes] Alongside the paper, we are publishing an open-source library of extended
 438 environments. This will be included in the supplemental material.
- 439 (d) Did you discuss whether and how consent was obtained from people whose data you’re
 440 using/curating? [N/A]
- 441 (e) Did you discuss whether the data you are using/curating contains personally identifiable
 442 information or offensive content? [N/A]
- 443 5. If you used crowdsourcing or conducted research with human subjects...
- 444 (a) Did you include the full text of instructions given to participants and screenshots, if
 445 applicable? [N/A]
- 446 (b) Did you describe any potential participant risks, with links to Institutional Review
 447 Board (IRB) approvals, if applicable? [N/A]
- 448 (c) Did you include the estimated hourly wage paid to participants and the total amount
 449 spent on participant compensation? [N/A]