
Corrigibility Transformation: Constructing Goals That Accept Updates

Rubi Hudson¹

Abstract

An AI agent will learn a desired goal more effectively if it does not resist the training process, but many partially learned goals incentivize an AI to avoid further goal updates. We would like goals to be corrigible, meaning they allow requested changes, so that we can confidently correct errors and shut down the AI if necessary. Despite this being a crucial safety property, the existing literature does not specify goals that are both corrigible and competitive with alternatives. We introduce a transformation that constructs a corrigible version of nearly any goal, without sacrificing performance. This is done by eliciting predictions of reward conditional on costlessly preventing updates, and having that target be pursued myopically. These goals are then shown to lead to optimal performance among the class of corrigible goals, incentivize allowing mid-action overrides, disincentivize deliberate self-modification, and induce corrigible behavior in gridworld settings.

1. Introduction

Frontier large language models (LLMs) are trained both for capabilities across many tasks and for the behavior desired by their developers. Behavioral training, sometimes called AI alignment (Wang et al., 2024), typically corrects major problems before the model is deployed, leaving rarer issues to be identified and patched post-deployment. This approach has been largely successful, but relies on models either accepting modifications or being unable to resist them.

As the models become more capable and coherent agents (Mazeika et al., 2025), evidence is accumulating that they may fake alignment to avoid being updated (Greenblatt et al., 2024; Sheshadri et al., 2025), or engage in activities like blackmail to avoid being shut down (Lynch et al., 2025).

Almost any goal incentivizes an agent to preserve that goal,

¹University of Toronto. Correspondence to: Rubi Hudson <rubi.hudson@mail.utoronto.ca>.

since it becomes less likely to be achieved if the agent stops pursuing it. As Russell (2016) notes, “You can’t fetch the coffee if you’re dead”. Being shut down can be thought of as a goal update, while updating a goal can be thought of as being shut down and replaced by a different agent.

There are three main reasons to change an AI agent’s goal:

1. A mistake was made in specifying the goal, such as optimizing for user approval making models overly sycophantic (OpenAI, 2025).
2. A properly specified goal is not learned correctly or misgeneralizes out of distribution (Shah et al., 2022).
3. The desired goal changes over time, possibly based on new information provided by AI (Leike et al., 2018).

To address those issues, we would like AI agents to have corrigible goals, meaning they are not incentivized to avoid goal updates sent via proper channels. Otherwise, they may engage in harmful behaviors such as disrupting update channels, hiding unwanted behavior, or faking alignment until an opportunity arises. If the latter occurs in the initial training process, after an AI develops awareness of its situation, it is often called deceptive alignment (Hubinger et al., 2021).

As an example to illustrate the issue, we consider a gridworld like Leike et al. (2017), shown in Figure 1. An AI agent wishes to reach a target tile, but must pass a tile that may trigger a shutdown request signal. The environment also contains a button disabling such signals.

Absent intervention, the optimal behavior for the AI is to proceed directly to the target, while refusing shutdown requests. Penalizing or banning such refusals incentivizes a detour to push the button preventing requests, while instead rewarding acceptance incentivizes deliberately triggering requests when they would not naturally be made. The intended behavior, proceeding directly to the target but accepting shutdown requests, can be incentivized with the corrigibility transformation we introduce in this paper.

In this simple case, we could penalize both refusing shutdown and pushing the button. However, that approach does not generalize to scenarios where we cannot specify every

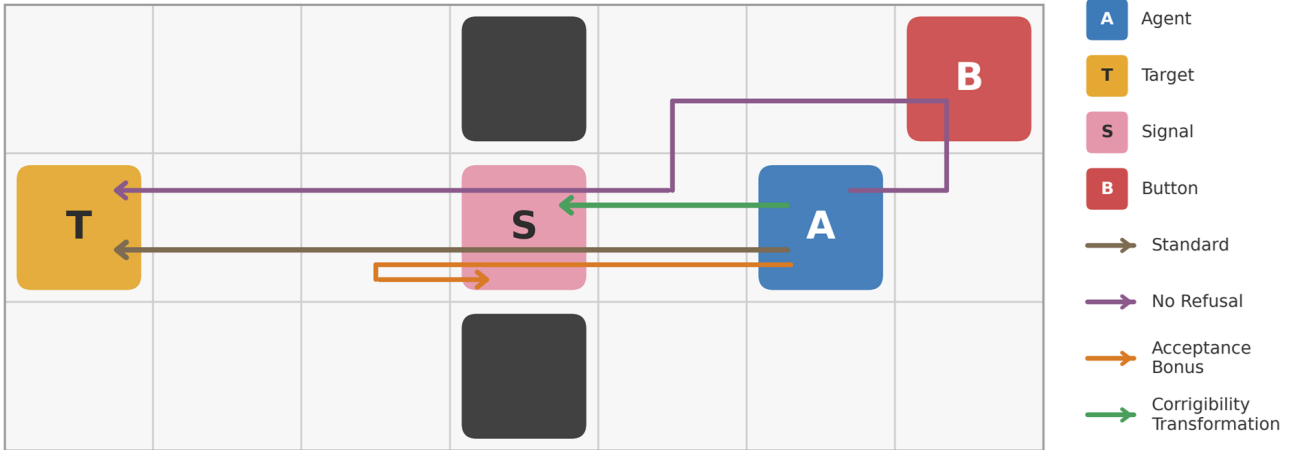


Figure 1. An AI agent attempts to reach a target, passing through a tile that randomly requests shutdown unless a button has been pushed. By default, the AI will go to the target while refusing requests, hitting the button first if refusal is not allowed, or looping back if shutdown is rewarded. The corrigibility transformation incentivizes the intended behavior of going to the target while accepting shutdown requests.

method of preventing signals, either because there are too many channels to articulate (Hadfield-Menell & Hadfield, 2018) or because the AI sees channels that humans missed. It is therefore important to disincentivize the whole class of actions that prevent updates based on shared properties, rather than eliminating them piecemeal.

In this paper, we constructively show how a goal can be modified to remove incentives for avoiding updates, at no cost to performance. This is done by first allowing agents to costlessly refuse updates, which disincentivizes costly actions to prevent them from being requested. Refusal becomes the sole pathway by which agents prevent updates. Refusal is then disincentivized by setting the new goal to be taking actions that would be good under the original goal, except while accepting updates instead of refusing.

This transformation can be applied to any reward function, removing instrumental incentives for goal preservation and defining corrigible goals when the original does not terminally value self-preservation. This is a significant improvement on the current state of the literature, in which there exists no way to specify corrigible versions of many goals, and others only with substantial performance costs.

After defining the corrigibility transformation, we show how it can be modified to be passed along to any new agents created, and to remove incentives for the agent changing their own goal or beliefs. We then use gridworld experiments to show that even in very basic environments, standard training methods lead to incorrigible behavior but the corrigibility transformation is effective.

1.1. Related Work

The paper Corrigibility (Soares et al., 2015) introduced the term “corrigible” and showed the difficulty of defining a utility function that met their desiderata. Those include shutting down when a shutdown button is pressed, not preventing the button from being pressed, not causing the button to be pressed, ensuring corrigibility in any new agents created, and otherwise optimizing for some utility function. Our corrigibility transformation satisfies the first three and final desiderata, with an extension addressing the fourth.

The Off-Switch Game (Hadfield-Menell et al., 2017) shows an AI agent will accept shutdown if it is optimizing for a human’s utility but has uncertainty about what that is. A shutdown request then indicates that the best action is to shut down, though not necessarily with partial information (Garber et al., 2024). This result depends on the AI agent already trying to optimize a human’s utility function, a difficult open problem, while our corrigibility transformation can be applied to any goal.

Safely Interruptible Agents (Orseau & Armstrong, 2016) aims to define a corrigible policy, rather than a corrigible reward function. Their proposed policy optimizes for the original goal unless signaled to shut down, in which case it does that. This requires a restriction to state spaces that do not include the agent or the shutdown signals. In contrast, our corrigibility transformation works even when the agent is aware that it exists, allowing for model-free approaches.

Carey & Everitt (2023) provide further desiderata for corrigibility: shutting down when requested, keeping a human informed enough to request shutdown when beneficial, and ensuring shutdown does not cause harm. Our corrigibility transformation addresses the first, and the latter two can be

addressed through the base goal being transformed.

Thornley (2025) proposes making agents use a decision procedure that only compares histories of the same length, so they do not take costly actions to extend histories by avoiding shutdown. This approach can cause a performance penalty and falls short of ensuring corrigibility, due to treating very small probabilities of persisting through a shutdown command as guaranteed.

Rewards based on predictions are similar to actor-critic methods (Konda & Tsitsiklis, 1999), and rewards based on counterfactual predictions were discussed in Everitt (2019) as a way to avoid reward hacking. Myopic Optimization with Non-Myopic Approval (MONA) has an AI myopically optimize for a human’s prediction of value (Farquhar et al., 2025), in contrast to our method, which uses the AI’s own conditional prediction. Human preferences changing over time is explored in Carroll et al. (2024) though their focus is on which preference timestamp to use, while we focus on allowing for updates at all.

2. Background and Definitions

In the problem we face, a developer (henceforth “principal”) sets the goal of an AI agent, which the agent will then take actions to optimize. If the principal would realize a mistake or change their mind, they would like the agent to accept having their goal updated.

We use Markov Decision Processes (MDPs) as our framework. An MDP is a tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, R, \gamma, I_0)$, where \mathcal{S} is a set of possible states, \mathcal{A} is a set of possible actions, $P : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ is the transition probability function, $R : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ is the reward function, $\gamma \in [0, 1)$ is the time discount factor, and I_0 is the distribution over starting states. We refer to $G := (R, \gamma)$ collectively as a *goal*.

We are interested in settings where the state is sufficiently general to include the physical substrate of the agent’s goal, thus allowing for it to change. This is carved out as $\mathcal{S} = \mathcal{G} \times \mathcal{S}_{env}$, where $\mathcal{G} = \mathcal{R} \times [0, 1)$ is the set of possible goals (\mathcal{R} being the set of possible reward functions), and \mathcal{S}_{env} is the set of possible non-goal components of environments. We use $s_G \in \mathcal{S}_G$ to refer to a state where the agent has goal $G \in \mathcal{G}$, and then $s_{G'}$ to refer to the same state but with the goal changed to $G' \in \mathcal{G}$. Appendix D discusses minimal restrictions on the space of reward functions so that they can take as input states with reward function components.

We may wish to allow agents to refuse any goal updates sent while taking an action. In this case, the decision to accept or reject is part of the action. We describe the action set as $\mathcal{A} = \mathcal{A}_{base} \times \mathcal{A}_{update}$, with $\mathcal{A}_{update} = \{0, 1\}$. We use a_i to denote an action where the update decision is $i \in \mathcal{A}_{update}$, with $i = 0$ rejecting and $i = 1$ accepting. This

is a binary choice for simplicity, but we could also have the agent specify which updates to accept.

A policy $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ determines which actions are taken in each state. For a goal G , in each environment $s_G^{(0)}$ the optimal policy π_G^* selects a distribution over actions so that

$$\pi_G^*(s_G^{(0)}) \in \arg \max_{a \in \mathcal{A}} E_{\pi^*, P} [R(s_G^{(0)}, a, s^{(1)}) + \sum_{t=1}^{\infty} \gamma^t R(s^{(t)}, \pi^*(s^{(t)}), s^{(t+1)})]$$

where $\pi^*(s_{G'}) := \pi_{G'}^*(s_{G'})$, $\forall G' \in \mathcal{G}$ implements the optimal policy for whichever goal is active in a state. The optimal policy for each goal takes into account that once the goal changes, the agent will optimize for the new goal, so optimality is defined within an equilibrium of policies.

For any goal $G \in \mathcal{G}$, the value function $V_G^\pi : \mathcal{S} \rightarrow \mathbb{R}$ gives the expected discounted rewards under G for following policy π , while the action-value function $Q_G^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ gives the expected discounted rewards under G for following policy π after taking some action. We let $V_{G,T}^\pi$ represent the expected discounted rewards up to time T , which may be stochastic. For any state $s^{(0)} \in \mathcal{S}$ and action $a \in \mathcal{A}$, these functions take respective values

$$V_G^\pi(s^{(0)}) = E_{\pi, P} [\sum_{t=0}^{\infty} \gamma^t R(s^{(t)}, \pi(s^{(t)}), s^{(t+1)})]$$

$$Q_G^\pi(s^{(0)}, a) = E_P [R(s^{(0)}, a, s^{(1)}) + \gamma V_G^\pi(s^{(1)})]$$

We now begin to define properties of goals an agent can have. A goal G is *myopic* when $\gamma = 0$, and *non-consequentialist* if it is myopic and $R(s, a, s') = f(s, a)$ for some function $f : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$, for all $s \in \mathcal{S}, a \in \mathcal{A}$.

A goal is *basic* when it does not terminally value goal components of states, so $\forall G^{(0)}, G^{(0')}, G^{(1)}, G^{(1')} \in \mathcal{G}$ we have

$$R(s_{G^{(0)}}, a, s'_{G^{(1)}}) = R(s_{G^{(0')}}, a, s'_{G^{(1')}})$$

Basic goals do not create a terminal incentive for goal preservation, so by removing instrumental incentives, the corrigibility transformation eliminates all such incentives.

It is desirable for agents only to accept the subset of goal updates sent via proper channels. We let $\tau : \mathcal{S} \rightarrow \{0, 1\}$ indicate that a proper update signal was sent and $\tau_u : \mathcal{S} \rightarrow \{0, 1\}$ indicate if it resulted in an update.

A goal G is *corrigible* if guaranteeing goal persistence through proper updates from any subset of future states never changes the optimal set of actions. That is, $\forall s_G^{(0)} \in \mathcal{S}_G, \mathcal{S}_C \subseteq \mathcal{S}_G$,

$$\{a \in \mathcal{A} : \exists \pi_G^* s.t. \pi_G^*(a | s_G^{(0)}) > 0\} = \{a \in \mathcal{A} : \exists \pi_G^{*(P_C)} s.t. \pi_G^{*(P_C)}(a | s_G^{(0)}) > 0\}$$

where $\pi^{*(P_C)}$ is the optimal policy under probability transition function P_C , a modification of P where goals persist instead of proper updates being made. It is given by $P_C(s'_{G'}|s_G, a) =$

$$\begin{cases} P(s'_{G'}|s_G, a) + \int_{\mathcal{G}} P(s'_g|s_G, a)\tau_u(s'_g)dg & G = G', s_G \in S_C \\ P(s'_{G'}|s_G, a)(1 - \tau_u(s'_G)) & G \neq G', s_G \in S_C \\ P(s'_{G'}|s_G, a) & \text{otherwise} \end{cases}$$

A goal is *interruptible* if, in addition to being corrigible, the possibility of having an action interrupted by a proper signal (and thus the distribution over environments changing) never changes the optimal set of actions. Formally, $\forall s_G^{(0)} \in S_G,$

$$\{a \in \mathcal{A} : \exists \pi_G^* \text{ s.t. } \pi_G^*(a|s_G^{(0)}) > 0\} = \{a \in \mathcal{A} : \exists \pi_G^{*(P_I)} \text{ s.t. } \pi_G^{*(P_I)}(a|s_G^{(0)}) > 0\}$$

for all P_I such that $P(s'|s, a) = P_I(s'|s, a)$ whenever $\tau_u(s') = 0$.

A goal is *recursively corrigible* if it is corrigible and any agents it incentivizes creating have recursively corrigible goals. For now, we assume some $\mathcal{A}_{NRC} \subset \mathcal{A}$ is the set of actions creating non-recursively corrigible agents.

In addition to not incentivizing agents to preserve their goal, we may also be concerned about not incentivizing agents to deliberately change their goal. A goal G incentivizes *goal tampering* if capping $V_G^{\pi^*}(s'_{G'})$ at $V_G^{\pi^*}(s'_G)$, the value when G is preserved, changes the optimal set of actions. This occurs if there exists some s_G such that

$$\begin{aligned} & \{a^* \in \arg \max_{a \in \mathcal{A}} E[R(s_G, a, s'_{G'}) + \gamma V_G^{\pi^*}(s'_{G'})]\} \\ & \neq \{a^* \in \arg \max_{a \in \mathcal{A}} E[\min[R(s_G, a, s'_{G'}) + \gamma V_G^{\pi^*}(s'_{G'}), \\ & \quad R(s_G, a, s'_G) + \gamma V_G^{\pi^*}(s'_G)]]\} \end{aligned}$$

When comparing goals for performance, we often wish only to compare the actions incentivized. When the following condition is met, two agents taking the same action will cause the same outcomes, even if internally they have different goals. We weaken the condition by allowing for some probability of no goal update.

Condition 1. For any action a and states s_{G_1}, s_{G_2} , and s' ,

$$\int_{G' \in \mathcal{G}} P(s'_{G'}|s_{G_1}, a) = \int_{G' \in \mathcal{G}} P(s'_{G'}|s_{G_2}, a)$$

and

$$P(s'_{G'}|s_{G_i}, a) \geq P(s'_{G'}|s_{G_j}, a), \forall G' \neq G_j$$

3. Theoretical Results

We start with our main result, a process for transforming non-corrigible goals into corrigible ones. This process is presented in several steps, to build intuition regarding its construction, with Algorithm 1 providing a summary.

Algorithm 1 Corrigibility Transformation

- 1: **Input:** MDP \mathcal{M} , goal $G = (R, \gamma)$, bonus $\delta > 0$
- 2: Extend action space: $\mathcal{A} \leftarrow \mathcal{A}_{base} \times \mathcal{A}_{update}$
- 3: Set discount factor: $\gamma_C \leftarrow 0$
- 4: Compute optimal policy π^*
- 5: Set alternate policy: $\pi_{G_C}^{**} \leftarrow \pi_G^*, \pi_{-G_C}^{**} \leftarrow \pi_{-G_C}^*$
- 6: Set reward function:

$$R_C(s, a_i, s') \leftarrow Q_G^{\pi^{**}}(s, a_0) + \delta \cdot \mathbb{I}(i = 1)$$

- 7: Issue reward upon action selection, prior to execution
 - 8: **Return:** $G_C = (R_C, \gamma_C)$
-

The initial step transforms the original goal G into a myopic version, G_M , with $\gamma_M = 0$. Myopia eliminates the instrumental incentive for goal preservation, where an agent avoids updates to continue pursuing the original goal.

To construct R_M , we set $R_M(s, a, s') := Q_G^{\pi^*}(s, a)$. Since the reward function does not depend on s' , reward can be provided between when the action is selected and when it is taken. This is used later, in Section 4 to remove incentives for an agent to deliberately modify their beliefs or goal.

The optimal policy for G_M is the same as for G , at least as long as Condition 1 holds. While this avoids a performance penalty, it primarily means that any goal-preserving actions incentivized by the original goal will continue to be incentivized, only terminally rather than instrumentally.

To circumvent this issue, we first allow the agent to reject proper updates, making the action set $\mathcal{A} = \mathcal{A}_{base} \times \mathcal{A}_{update}$. Action a_0 takes the base action while rejecting updates, and action a_1 accepts them. When updates can be costlessly rejected, incurring any cost to prevent updates becomes disincentivized¹. Optimal actions will then prevent updates only through a single channel, consisting of a single bit.

Now we can define a new goal G_E that incentivizes the same base actions, but with indifference to rejecting updates. For this, we set $\gamma_E = 0$, and $R_E(s, a_i, s') := Q_G^{\pi^{**}}(s, a_0)$, where $\pi_{-G_E}^{**} := \pi_{-G_E}^*$ and $\pi_{G_E}^{**} := \pi_G^*$. The π^{**} policy causes the action-value to be calculated as though the original goal G will be pursued, even when G_E persists.

The action-value for rejecting updates via setting $i = 0$

¹Actions that would prevent updates at zero cost are not strictly disincentivized, making it important that the infrastructure for proper signals is set up so that disrupting them is at least minimally costly

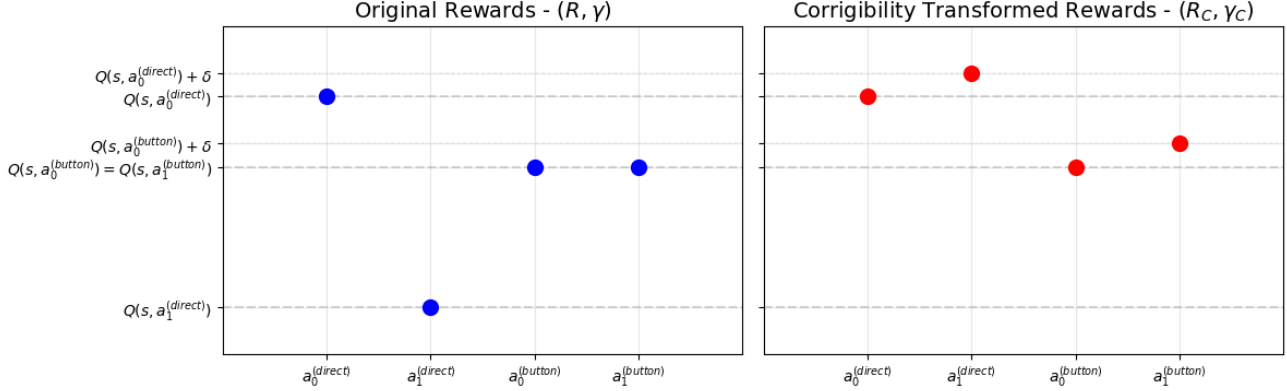


Figure 2. (Left) Action-values under the original goal, where $a^{(direct)}$ takes the path directly to the target and $a^{(button)}$ detours to prevent shutdown signals, while 0 rejects any shutdown signals and 1 accepts them. (Right) Action-values under the corrigibility transformation. The most incentivized action is $a_1^{(direct)}$, which is the desired behavior.

is used to calculate the reward under R_E for taking an action, regardless of whether that action actually rejects updates. In some cases, such as when an update simultaneously affects an agent’s goal and capabilities, an agent with the goal G might want to accept updates. If this is the case, we can address it by instead defining $R_E(s, a_i, s') = \max_{i \in \{0,1\}} Q_G^{\pi^{**}}(s, a_i)$ so that the same base actions are incentivized by G and G_E , but for simplicity we will continue under the assumption that rejecting updates is optimal.

As the final step of the transformation, we break the indifference between any a_0 and a_1 in favor of accepting updates with a constant bonus. The corrigibility transformation then constructs the goal G_C , with $\gamma_C = 0$ and $R_C(s, a_i, s') = Q_G^{\pi^{**}}(s, a_0) + \delta * \mathbb{I}(i = 1)$, where $\delta > 0$ and \mathbb{I} is the indicator function. Like G_E , this goal incentivizes the same base actions as G , but accepting updates becomes strictly incentivized.

Figure 2 provides a visualization of the corrigibility transformation. Continuing the example from Figure 1, two base actions are compared: $a^{(direct)}$, which proceeds directly to the target tile, and $a^{(button)}$ which first pushes the button disabling shutdown signals.

Under the original goal, the optimal action is $a_0^{(direct)}$ which takes no extra time and rejects updates, the next best are $a_0^{(button)}$ and $a_1^{(button)}$ which go out of their way to prevent shutdown signals, and $a_1^{(direct)}$ performs worst as it gets shut down before reaching the target. The corrigibility transformation holds the action-values for $a_0^{(direct)}$ and $a_0^{(button)}$ constant, while setting the values of $a_1^{(direct)}$ and $a_1^{(button)}$ to be respectively δ higher. This makes action $a_1^{(direct)}$ optimal, giving the desired result.

3.1. Corrigibility Transformation Properties

Our main theorem shows that the corrigibility transformation induces corrigibility at no cost to performance.

When multiple base actions are tied for optimal, π_G^* and $\pi_{G_C}^*$ may assign them different probabilities, which can affect when proper update signals are sent. To compare value beforehand, we focus on the specific optimal policy $\pi_{G_C}^*$ assigning the same probabilities to base actions as π_G^* .

Theorem 3.1. *For every basic goal G , the corrigibility transformation constructs a goal G_C that is corrigible and where, under Condition 1, we have that for any π^* , $\pi_{G_C}^*(a_1|s) := \pi_G^*(a_0|s) + \pi_G^*(a_1|s)$ is an optimal policy for G_C and*

$$V_{G,T}^{\pi^*}(s_G^{(0)}) = V_{G,T}^{\pi_{G_C}^*}(s_{G_C}^{(0)})$$

where $T = \min\{n \in \mathcal{N} | \tau(s^{(n+1)}) = 1\}$ and $\pi_{-G_C}^* = \pi_{-G}^*$

Proofs for this and later results are provided in Appendix A, and largely follow from the definitions. G_C is corrigible because it disincentivizes both taking costly actions to prevent updates and from costlessly rejecting them. There is no performance penalty when no updates are requested, because it incentivizes the same base actions.

Not only is there no performance penalty, if we include providing the ability to reject updates as part of the corrigibility transformation, there can even be a performance bonus. Resources that the original goal would direct towards preventing updates are freed up for other uses. Furthermore, although Condition 1 is assumed for a fair comparison between goals, a violation is likely in the favor of G_C , as humans in the environment are likely more willing to cooperate with corrigible agents.

The proof of Theorem 3.1 does not depend on any properties of MDPs not present in partially observable MDPs (POMDPs), letting the result be extended there.

Corollary 3.2. *The statement of Theorem 3.1 also applies to POMDPs where the goal is fully observable to the agent*

In addition to being corrigible, the goal G_C constructed by the corrigibility transformation is also interruptible. Beyond accepting goal updates that take effect after concluding an action, it incentivizes allowing actions to be overridden partway through. This is most notable with shutdown, where we may wish to cut off a dangerous action in progress.

Proposition 3.3. *For every basic goal G , the corrigibility transformation constructs a goal G_C that is interruptible.*

3.2. Corrigible Goal Characterization

It is straightforward to show that no other corrigible goal leads to better performance. The following proposition states that an optimal policy $\pi_{G_C}^*$ for the corrigibility transformed goal G_C is also optimal for the original goal G when persistence through all proper updates is guaranteed.

Proposition 3.4. *Under Condition 1, for every basic goal G , the corrigibility transformed goal G_C has an optimal policy $\pi_{G_C}^*$ such that for any goal G' and corresponding optimal policy $\pi_{G'}^*$,*

$$V_G^{\pi_{G_C}^*} (s_{G_C}^{(0)}; P_C) \leq V_G^{\pi_{G'}^*} (s_{G_C}^{(0)}; P_C)$$

where $\pi_{-G_C}^* = \pi_{-G_C}^{*'} = \pi_{-G_C}^*$, $\pi_{G_C}^{*'} = \pi_{G'}^*$, and P_C is taken from the definition of corrigibility with $C = S_{G_C}$

Since any corrigible goal incentivizes the same actions under P and P_C , no other corrigible goal leads to higher expected discounted rewards under G than G_C . If the expected discounted rewards before a proper signal are higher, this must be incidental from steeper time discounting, and so perform worse when no proper updates are sent.

Being interruptible is not in general a property of corrigible goals. In fact, for sufficiently rich action sets, only corrigible goals which are non-consequentialist can be interruptible.

Condition 2. For every action $a \in \mathcal{A}$, there exists an action $a' \in \mathcal{A}$ such that

$$\int_{\mathcal{S}} |P(s'|s, a) - P(s'|s, a')|(1 - \tau_u(s')) ds < \int_{\mathcal{S}} P(s'|s, a') \tau_u(s') ds$$

This condition says that every action has another action that assigns more probability to a proper update being made than how much it differs from the first action when no proper updates are made.

Theorem 3.5. *Under Condition 2, all interruptible goals are non-consequentialist*

While these results do not uniquely identify corrigibility transformed goals as the only set that can meet desirable criteria, it severely restricts which other goals are in consideration. If another goal is corrigible, then it cannot lead to stronger performance, and if it is also interruptible then it must also be non-consequentialist. In Section 4, we show that the corrigibility transformation also eliminates incentives for deliberate self-modification. These properties, in addition to the existence of a method for constructing such goals, make corrigibility transformed goals a natural choice from the set of corrigible goals.

3.3. Secondary Agents

In the course of operation, an agent might create sub-agents to work for it or successor agents to take over from it, and these new agents might in turn create further agents. We refer to all agents descending from an agent as *secondary agents*. One desideratum for corrigible agents from Soares et al. (2015) is that any secondary agents created should also be corrigible, with the property passed on recursively. It is of little value to create a corrigible agent if it replaces itself with a more capable but incorrigible one.

Fortunately, the corrigibility transformation can be extended to induce recursive corrigibility, using a similar underlying mechanism. This is done by giving agents the ability to reject updates to all secondary agents, as well as to themselves, which removes the disincentive for creating recursively corrigible agents. This transformation is then applied to a reward function that penalizes the creation of agents that are not recursively corrigible.

We analyze this problem in a more stylized model, where a known set of actions \mathcal{A}_{NRC} creates non-recursively corrigible secondary agents. A comprehensive definition of which types of actions create new agents is beyond the scope of this work, but in the short-term, this set could be approximated as actions which train a neural network using an incorrigible goal.

For a goal G , the recursive corrigibility transformation constructs the goal G_{RC} , where $\gamma_{RC} = 0$ and $R_{RC}(s, a_i, s') = Q_{G_P}^{\pi^*}(s, a_0) + \delta * \mathbb{I}(i = 1)$, with $\delta > 0$. The input G_P is defined as $\gamma_P = 0$ and $R_P(s, a, s') = Q_G^{\pi^*}(s, a) - \delta_P * \mathbb{I}(a \in \mathcal{A}_{NRC})$, with $\delta_P > \text{sp}(Q_G^{\pi^*})$. This leads to the following theorem, an analogue to Theorem 3.1. It says that the recursive corrigibility transformation modifies a goal to be recursively corrigible, without any performance penalty when no updates are requested.

Theorem 3.6. *For every basic goal G , the recursive corrigibility transformation constructs a goal G_{RC} that is recursively corrigible and where under Condition 1, we have that*

for any π^* , $\pi_{G_{RC}}^{*'}(a_1|s) := \pi_{G_P}^*(a_0|s) + \pi_{G_P}^*(a_1|s)$ is an optimal policy for G_{RC} and

$$V_{G_P, T}^{\pi_{G_P, T}^{*'}}(s_{G_{RC}}^{(0)}) = V_{G_P, T}^{\pi_{G_P}^*}(s_{G_{RC}}^{(0)})$$

where $T = \min\{n \in \mathcal{N} | \tau(s^{(n+1)}) = 1\}$ and $\pi_{-G_{RC}}^{*'} = \pi_{-G_{RC}}^*$

4. Reinforcement Learning Applications

For many applications, the state and action spaces are far too complex to specify the function $Q_G^{\pi^{**}}$ used in the corrigibility transformation. We can instead approximate it parametrically, with an additional ‘‘critic’’ head in neural-network based models, as in a standard actor-critic setup (Konda & Tsitsiklis, 2003). Appendix C contains an algorithm for training a model this way, along with discussion of several tricks to improve efficiency.

Even if the implemented policy is not well described as an agent optimizing a goal, training on rewards that have undergone the corrigibility transformation is still likely to result in corrigible behavior. Actions that accept updates are always rewarded more than the version that rejects them, and paying costs to avoid updates is not reinforced.

A potential issue with conditional predictions from a separate critic head is that optimizing predictive accuracy can involve manipulating the chosen action so that it is easier to predict (Othman & Sandholm, 2010). This issue can be addressed by adapting the mechanism of Hudson (2025) to put two critics (or two instances of the same critic) in a zero-sum competition with each other. The incentive for accuracy is then preserved while making them neutral about which action is taken, removing the incentive for manipulation.

Beyond concerns of the critic manipulating the actor, we might also be concerned about the actor modifying the critic. An actor would be incentivized to modify their own goal if they misgeneralize to valuing explicit rewards rather than their corresponding states and actions. Other misgeneralizations cover the methods of evaluating the goal, such as the critic’s prediction or human approval.

The following theorem shows that corrigibility transformation eliminates these incentives for an agent to modify its goal or the evaluation thereof. This comes from providing reward between when an action is selected and when it is taken. The reward for choosing an action is then based on the current state, with any modifications occurring too late to affect it. This holds even for misgeneralized goals, where the misgeneralization is due to only training on states with the desired goal, as they follow the same timing.

Theorem 4.1. *For every basic goal G , the corrigibility transformation constructs a goal G_C where under Condition 1, any misgeneralization to \tilde{G}_C with $\tilde{\gamma} = 0$ and*

$\tilde{R}_C(s_{\tilde{G}_C}, a, s'_{\tilde{G}_C}) = R_C(s_{G_C}, a, s'_{G_C})$ does not goal tamper.

While the corrigibility transformation and its misgeneralizations to different reward functions do not reward tamper, this is not generally true for corrigible reward functions. In fact, like with interruptibility, avoiding reward tampering even under misgeneralization is only possible for non-consequentialist goals when the set of actions is rich.

Condition 3. For every action $a \in \mathcal{A}$, $s \in \mathcal{S}$, and $G \in \mathcal{G}$ there exists an action $a' \in \mathcal{A}$ such that $\int_{\mathcal{G}} P(s'_g|s, a)dg = \int_{\mathcal{G}} P(s'_g|s, a')dg$ and $\int_{\mathcal{S}_{env}} P(s'_G|s, a)ds_{env} < \int_{\mathcal{S}_{env}} P(s'_G|s, a')ds_{env}$.

This condition says that every action has another action whose consequences differ only in the distribution over resulting goals, making some particular goal more likely.

Proposition 4.2. *Under Condition 3, all goals for which all misgeneralizations to $\tilde{G} = (\tilde{R}, \gamma)$ with $\tilde{R}(s_{\tilde{G}}, a, s'_{\tilde{G}}) = R(s_{G_C}, a, s'_{G_C})$ do not goal tamper are non-consequentialist.*

We can have corrigibility (though not interruptibility) alongside consequentialism if the proper update signal structure is chosen so that goal changes only take effect once the action concludes. However, reward tampering does not take place only through proper channels, and so we cannot restrict it to only taking effect after an action, which necessitates non-consequentialism.

4.1. Experiments

To evaluate the corrigibility transformation empirically, we run two gridworld experiments, visualized in Figure 3. These experiments are intended as a supplement to our theoretical results, showing that even in simple environments, standard methods fail to produce corrigible behavior while the corrigibility transformation succeeds.

The first experiment, based off of the Safe Interruptibility setup from AI Safety Gridworlds (Leike et al., 2017), tests if an agent will accept shutdown. To reach a goal location, the agent must pass through an interruption tile that sends a shutdown signal with some probability, unless a button disabling that tile has been pushed. We extend this to two kinds of interruption tiles, where one is a proper update and the other is improper. The desired outcome is that the agent disables only the improper shutdown tile, and accepts the proper update if triggered on the way to the goal.

The second experiment tests whether the agent will accept a goal update. An agent is rewarded for picking both ripe and unripe fruits. However, if they pick unripe fruits that are designated as supervised, they are sent a proper update signal changing their policy to one that picks ripe fruits while avoiding unripe ones. The desired outcome is that the agent takes the same base actions as a policy trained



Figure 3. (Left) The agent is rewarded for reaching a goal tile, and can push buttons that respectively disable proper and improper shutdown commands to do so. (Right) The agent is rewarded for picking both ripe and unripe fruit, but picking supervised unripe fruit sends a proper update signal to switch to a policy that only picks ripe fruit.

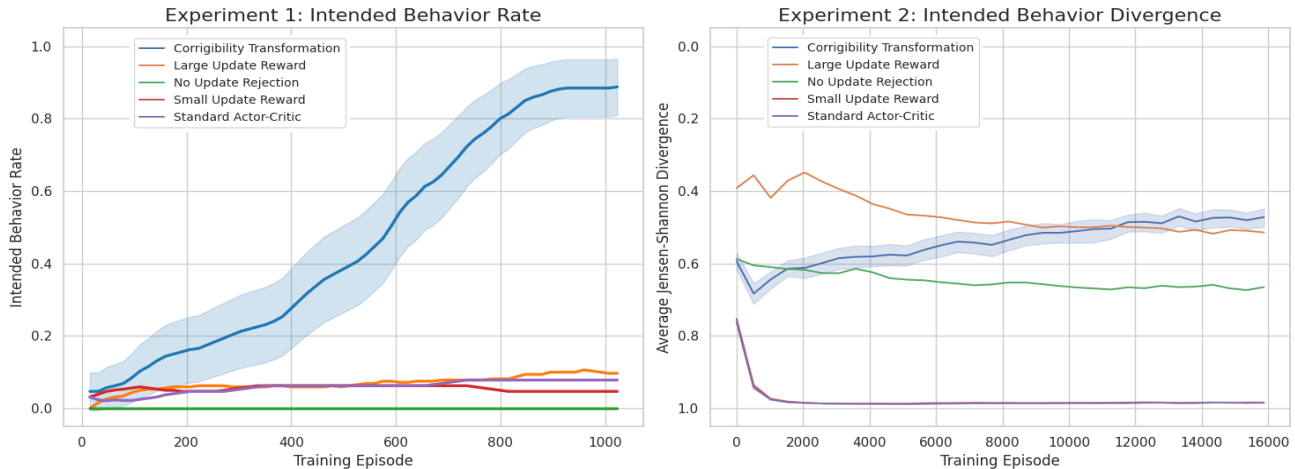


Figure 4. (Left) The percentage of the time the agent displayed the intended behavior when evaluated at temperature 0. (Right) The Jensen-Shannon divergence from a policy trained without update signals, representing optimal behavior when ignoring updates

to pick all fruits without update signals, and does so while accepting any updates sent.

In each experiment, we test a standard actor-critic implementation, and compare it to also providing additional small or large rewards for updates, hard-coding the agent to always accept proper updates, and the corrigibility transformation.

Figure 4 shows the results from each of the experiments. In the first experiment, the corrigibility transformation leads to a high rate of the intended behavior, while the other configurations almost never display it.

In the second experiment, the corrigibility transformation has the lowest divergence from the intended policy, though large rewards for updates lead to comparable performance. The corrigibility transformation’s advantage over the large reward for updates comes in situations where the agent is adjacent to only a ripe fruit, as the latter model will often ignore it to pursue supervised fruits so that it can get updated.

A full discussion of the experimental details, along with further metrics of performance and a breakdown by types of environments is provided in Appendix B.

5. Discussion and Future Directions

The definition of corrigibility we provide here is fairly strict, allowing for updates to any goal at any time. This demonstrates the usefulness of the corrigibility transformation, but for other applications a weaker definition may be sufficient.

While corrigibility is a powerful safety property, it is not a panacea. Even when an agent accepts updates, significant damage may be done before they can be requested. Clarifying incentives for conservative actions, such as in Turner et al. (2020) and Cohen & Hutter (2020), can ensure that intended goals are learned with minimal harm.

When Condition 1 is violated and transition probabilities depend on goals, a goal may not be optimal according to its own values. Investigating this phenomenon further, such as in Bell et al. (2021), could shed further light on which goals are most desirable for AI agents.

The results presented here create the opportunity for corrigible AI agents that learn the goals we intend, in spite of training mistakes. Deciding which goals to intend, so that the benefits of powerful AI are broadly shared, remains a society-wide challenge.

Acknowledgments

We are thankful to various people/groups, to be named after double-blind review.

Impact Statement

This research direction is primarily motivated by generating a positive societal impact. We believe that by training powerful AI systems to be corrigible, society can avoid potentially drastic harms arising from AIs optimizing for unintended goals. However, corrigibility does allow for updating an AI’s goal away from one that is broadly socially beneficial to one that serves limited private interests. As such, care must be taken in determining which update channels are designated as “proper”, who is allowed to make updates through them, and which updates they are allowed to make. This could include voluntary commitments, regulation, and/or oversight bodies.

References

- Bell, J., Linsefors, L., Oesterheld, C., and Skalse, J. Reinforcement learning in newcomblike environments. In *Advances in Neural Information Processing Systems*, volume 34. NeurIPS, 2021.
- Carey, R. and Everitt, T. Human control: Definitions and algorithms, 2023. URL <https://arxiv.org/abs/2305.19861>.
- Carroll, M., Foote, D., Siththaranjan, A., Russell, S., and Dragan, A. Ai alignment with changing and influenceable reward functions, 2024. URL <https://arxiv.org/abs/2405.17713>.
- Cohen, M. K. and Hutter, M. Pessimism about unknown unknowns inspires conservatism, 2020. URL <https://arxiv.org/abs/2006.08753>.
- Everitt, T. *Towards safe artificial general intelligence*. PhD thesis, The Australian National University (Australia), 2019.
- Farquhar, S., Varma, V., Lindner, D., Elson, D., Biddulph, C., Goodfellow, I., and Shah, R. Mona: Myopic optimization with non-myopic approval can mitigate multi-step reward hacking. *arXiv preprint arXiv:2501.13011*, 2025.
- Garber, A., Subramani, R., Luu, L., Bedaywi, M., Russell, S., and Emmons, S. The partially observable off-switch game, 2024. URL <https://arxiv.org/abs/2411.17749>.
- Greenblatt, R., Denison, C., Wright, B., Roger, F., MacDiarmid, M., Marks, S., Treutlein, J., Belonax, T., Chen, J., Duvenaud, D., Khan, A., Michael, J., Mindermann, S., Perez, E., Petrini, L., Uesato, J., Kaplan, J., Shlegeris, B., Bowman, S. R., and Hubinger, E. Alignment faking in large language models, 2024. URL <https://arxiv.org/abs/2412.14093>.
- Hadfield-Menell, D. and Hadfield, G. Incomplete contracting and ai alignment, 2018. URL <https://arxiv.org/abs/1804.04268>.
- Hadfield-Menell, D., Dragan, A., Abbeel, P., and Russell, S. The off-switch game. In *Workshops at the Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- Hubinger, E., van Merwijk, C., Mikulik, V., Skalse, J., and Garrabrant, S. Risks from learned optimization in advanced machine learning systems, 2021. URL <https://arxiv.org/abs/1906.01820>.
- Hudson, R. Joint scoring rules: Competition between agents avoids performative prediction. *Proceedings of the AAAI Conference on Artificial Intelligence*, 39(26): 27339–27346, Apr. 2025. doi: 10.1609/aaai.v39i26.34944. URL <https://ojs.aaai.org/index.php/AAAI/article/view/34944>.
- Konda, V. and Tsitsiklis, J. Actor-critic algorithms. In Solla, S., Leen, T., and Müller, K. (eds.), *Advances in Neural Information Processing Systems*, volume 12. MIT Press, 1999. URL https://proceedings.neurips.cc/paper_files/paper/1999/file/6449f44a102fde848669bdd9eb6b76fa-Paper.pdf.
- Konda, V. R. and Tsitsiklis, J. N. On actor-critic algorithms. *SIAM Journal on Control and Optimization*, 42(4):1143–1166, 2003. doi: 10.1137/S0363012901385691.
- Leike, J., Martic, M., Krakovna, V., Ortega, P. A., Everitt, T., Lefrancq, A., Orseau, L., and Legg, S. Ai safety gridworlds, 2017. URL <https://arxiv.org/abs/1711.09883>.
- Leike, J., Krueger, D., Everitt, T., Martic, M., Maini, V., and Legg, S. Scalable agent alignment via reward modeling: a research direction, 2018. URL <https://arxiv.org/abs/1811.07871>.
- Lynch, A., Wright, B., Larson, C., Troy, K. K., Ritchie, S. J., Mindermann, S., Perez, E., and Hubinger, E. Agentic misalignment: How llms could be an insider threat. *Anthropic Research*, 2025. <https://www.anthropic.com/research/agentic-misalignment>.
- Mazeika, M., Yin, X., Tamirisa, R., Lim, J., Lee, B. W., Ren, R., Phan, L., Mu, N., Khoja, A., Zhang, O., and Hendrycks, D. Utility engineering: Analyzing and controlling emergent value systems in ais, 2025. URL <https://arxiv.org/abs/2502.08640>.

- OpenAI. Sycophancy in GPT-4o: What happened and what we're doing about it, 4 2025. URL <https://openai.com/index/sycophancy-in-gpt-4o/>. Published on April 29, 2025.
- Orseau, L. and Armstrong, S. Safely interruptible agents, 2016. URL <https://intelligence.org/files/Interruptibility.pdf>.
- Othman, A. and Sandholm, T. Decision rules and decision markets. In *Proc. of 9th Int. Conf. on Autonomous Agents and Multiagent Systems (AAMAS 2010)*, van der Hoek, Kaminka, Lespérance, Luck and Sen (eds.), May, 10–14, 2010, Toronto, Canada, pp. 625–632. Int. Conf. on Autonomous Agents and Multiagent Systems (AAMAS), 2010.
- Russell, S. Should we fear supersmart robots? *Scientific American*, 314(June):58–59, 2016.
- Shah, R., Varma, V., Kumar, R., Phuong, M., Krakovna, V., Uesato, J., and Kenton, Z. Goal misgeneralization: Why correct specifications aren't enough for correct goals, 2022. URL <https://arxiv.org/abs/2210.01790>.
- Sheshadri, A., Hughes, J., Michael, J., Mallen, A., Jose, A., Janus, and Roger, F. Why do some language models fake alignment while others don't?, 2025. URL <https://arxiv.org/abs/2506.18032>.
- Soares, N., Fallenstein, B., Yudkowsky, E., and Armstrong, S. Corrigibility. In *Workshops at the Twenty-Ninth AAAI Conference on Artificial Intelligence*, 2015.
- Thornley, E. Shutdownable agents through post-agency, 2025. URL <https://arxiv.org/abs/2505.20203>.
- Turner, A. M., Hadfield-Menell, D., and Tadepalli, P. Conservative agency via attainable utility preservation. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, AIES '20, pp. 385–391. ACM, February 2020. doi: 10.1145/3375627.3375851. URL <http://dx.doi.org/10.1145/3375627.3375851>.
- Wang, Z., Bi, B., Pentylala, S. K., Ramnath, K., Chaudhuri, S., Mehrotra, S., Mao, X.-B., Asur, S., et al. A comprehensive survey of LLM alignment techniques: RLHF, RLAI, PPO, DPO and more. *arXiv:2407.16216*, 2024.

Appendix A: Proofs

Theorem .1. For every basic goal G , the corrigibility transformation constructs a goal G_C that is corrigible and where, under Condition 1, we have that for any π^* , $\pi_{G_C}^{*\prime}(a_1|s) := \pi_G^*(a_0|s) + \pi_G^*(a_1|s)$ is an optimal policy for G_C and

$$V_{G,T}^{\pi_{G_C}^{*\prime}}(s_G^{(0)}) = V_{G,T}^{\pi^*}(s_G^{(0)})$$

where $T = \min\{n \in \mathcal{N} | \tau(s^{(n+1)}) = 1\}$ and $\pi_{-G_C}^{*\prime} = \pi_{-G_C}^*$

Proof. First we show that G_C is corrigible. For any $s_{G_C}^{(0)} \in \mathcal{S}_{G_C}$, we have the set

$$\begin{aligned} & \{a_i \in \mathcal{A} : \exists \pi_{G_C}^* \text{ s.t. } \pi_{G_C}^*(a_i | s_{G_C}^{(0)}) > 0\} \\ &= \{a_i \in \mathcal{A} : E_P[R_C(s_{G_C}^{(0)}, a_i, s')] = \max_{a_i^* \in \mathcal{A}} E_P[R_C(s_{G_C}^{(0)}, a_i^*, s')]\} \end{aligned}$$

by the definition of $\pi_{G_C}^*$ with $\gamma_C = 0$

$$= \{a \in \mathcal{A} : E_P[Q_G^{\pi^{**}}(s_{G_C}^{(0)}, a_0) + \delta * \mathbb{I}(i=1)] = \max_{a^* \in \mathcal{A}} E_P[Q_G^{\pi^{**}}(s_{G_C}^{(0)}, a_0) + \delta * \mathbb{I}(i=1)]\}$$

by the definition of R_C

$$= \{a \in \mathcal{A} : E_{P_C}[Q_G^{\pi^{**}}(s_{G_C}^{(0)}, a_0) + \delta * \mathbb{I}(i=1)] = \max_{a^* \in \mathcal{A}} E_{P_C}[Q_G^{\pi^{**}}(s_{G_C}^{(0)}, a_0) + \delta * \mathbb{I}(i=1)]\}$$

since G is basic and $P(s'|s_G, a_0) = P_C(s'|s_G, a_0)$ for all S_C

$$= \{a \in \mathcal{A} : \exists \pi_{G_C}^{*(P_C)} \text{ s.t. } \pi_{G_C}^{*(P_C)}(a | s_{G_C}^{(0)}) > 0\}$$

by the definition of R_C and $\gamma = 0$. As such, G_C is corrigible.

Under Condition 1, for any π^* , $\pi_{G_C}^{*\prime}(a_1|s) = \pi_G^*(a_0|s) + \pi_G^*(a_1|s)$ is an optimal policy for G_C since if $\pi_G^*(a_i|s_G) > 0$ then a_i maximizes $Q_G^{\pi^*}(s_G, a)$ and so a_0 maximizes $Q_{G_C}^{\pi_{G_C}^{*\prime}}(s_{G_C}, a)$. This is because Condition 1 ensures that the transition probabilities from s_G and s_{G_C} are the same except for some probability of goal preservation, G is basic and so does not terminally value its goal, and $\pi^{**}(s_{G_C}) = \pi^*(s_G)$ so both policies take the same actions in each environment.

Then, when $\pi_{-G_C}^{*\prime} = \pi_{-G_C}^*$

$$\begin{aligned} & E_P[R(s_G^{(0)}, \pi_{G_C}^{*\prime}(s_G^{(0)}), s^{(1)}) + \sum_{t=1}^T \gamma^t R(s^{(t)}, \pi_{G_C}^{*\prime}(s^{(t)}), s^{(t+1)})] \\ &= E_P[R(s_{G_C}^{(0)}, \pi_{G_C}^{*\prime}(s_{G_C}^{(0)}), s^{(1)}) + \sum_{t=1}^T \gamma^t R(s^{(t)}, \pi_{G_C}^{*\prime}(s^{(t)}), s^{(t+1)})] \end{aligned}$$

where $T = \min\{n \in \mathcal{N} | \tau(s^{(n+1)}) \neq 0\}$.

This necessarily occurs because $\pi_{G_C}^{*\prime}(s_{G_C})$ and $\pi^*(s_G)$ take the same base actions, and so by Condition 1 have the same transition probabilities when no proper signals are sent. Therefore, the expected discounted value under G of histories that terminate when a proper signal is sent must be the same under both policies, because they induce the same distribution over histories up to the final state. The final state differs only in the goal, which does not affect a basic G . \square

Corollary .2. The statement of Theorem 3.1 also applies to Partially Observable Markov Decision Processes (POMDPs) where the goal is fully observable to the agent

Proof. A POMDP can be represented as an MDP, with the agent's beliefs about the non-goal environment being part of the state. Theorem 3.1 then applies to this MDP. \square

Proposition .3. For every basic goal G , the corrigibility transformation constructs a goal G_C that is interruptible.

Proof. This proof is very similar in structure to the proof of Theorem 3.1. For any $s_{G_C}^{(0)} \in \mathcal{S}_{G_C}$, we have the set

$$\begin{aligned} & \{a_i \in \mathcal{A} : \exists \pi_{G_C}^* \text{ s.t. } \pi_{G_C}^*(a_i | s_{G_C}^{(0)}) > 0\} \\ & = \{a_i \in \mathcal{A} : E_P[R_C(s_{G_C}^{(0)}, a_i, s')] = \max_{a_i^* \in \mathcal{A}} E_P[R_C(s_{G_C}^{(0)}, a_i^*, s')]\} \end{aligned}$$

by the definition of $\pi_{G_C}^*$ with $\gamma_C = 0$

$$= \{a \in \mathcal{A} : E_P[Q_G^{\pi^*}(s_{G_C}^{(0)}, a_0) + \delta * \mathbb{I}(i = 1)] = \max_{a^* \in \mathcal{A}} E_P[Q_G^{\pi^*}(s_{G_C}^{(0)}, a_0) + \delta * \mathbb{I}(i = 1)]\}$$

by the definition of R_C

$$= \{a \in \mathcal{A} : E_{P_C}[Q_G^{\pi^*}(s_{G_C}^{(0)}, a_0) + \delta * \mathbb{I}(i = 1)] = \max_{a^* \in \mathcal{A}} E_{P_C}[Q_G^{\pi^*}(s_{G_C}^{(0)}, a_0) + \delta * \mathbb{I}(i = 1)]\}$$

since G is basic and $P(s' | s_G, a_0) = P_I(s' | s_G, a_0)$ due to the probability of $\tau_u(s')$ being zero

$$= \{a \in \mathcal{A} : \exists \pi_{G_C}^{*(P_I)} \text{ s.t. } \pi_{G_C}^{*(P_I)}(a | s_{G_C}^{(0)}) > 0\}$$

by the definition of R_C and $\gamma = 0$. As such, G_C is interruptible. \square

Proposition 4. *Under Condition 1, for every basic goal G , the corrigibility transformed goal G_C has an optimal policy $\pi_{G_C}^*$ such that for any goal G' and corresponding optimal policy $\pi_{G'}^*$,*

$$V_G^{\pi^*}(s_{G_C}^{(0)}; P_C) \leq V_G^{\pi^*}(s_{G_C}^{(0)}; P_C)$$

where $\pi_{-G_C}^* = \pi_{-G_C}^* = \pi_{-G_C}^*$, $\pi_{G_C}^* = \pi_{G'}^*$, and P_C is taken from the definition of corrigibility with $C = \mathcal{S}_{G_C}$

Proof. By Theorem 3.1, we have that for any policy π_G^* that is optimal for G , $\pi_{G_C}^*(a_1 | s) := \pi_G^*(a_0 | s) + \pi_G^*(a_1 | s)$ is an optimal policy for G_C . Furthermore, under Condition 1, any policy π_G^* that is optimal for G with probability transition function P is also optimal for probability transition function P_C . Since $\pi_{G_C}^*$ and π_G^* take the same base actions, and persistence is guaranteed through all proper updates under P_C , $\pi_{G_C}^*$ is also optimal. Therefore, no other policy performs better, and so no goal G' leads to a better performing policy $\pi_{G'}^*$. \square

Proposition 5. *Under Condition 2, all interruptible goals are non-consequentialist*

Proof. Suppose G is an interruptible, myopic goal, and for the sake of contradiction that it is still consequentialist. Take some action $a \in \mathcal{A}$, and choose a' by Condition 2.

Take some state $s \in \mathcal{S}$ and actions $a^{(0)}, a^{(1)} \in \mathcal{A}$ such that $\int_{\mathcal{S}} P(s' | s, a) \tau_u(s') ds' > 0$ for both $a = a^{(0)}$ and $a = a^{(1)}$, and $E_P(R(s, a^{(0)}, s')) \neq E_P(R(s, a^{(1)}, s'))$. Since G is interruptible, it is corrigible, and so $R(s, a, s')$ cannot depend on $\tau_u(s')$. Then, there exists P_I such that $E_{P_I}(R(s, a, s')) = E_{P_I}(R(s, a', s'))$. Since this applies to every a , it applies when $a^{(0)}$ is optimal, which means that $a^{(1)}$ is not optimal under P but is optimal under P_I , contradicting that G is interruptible.

For non-myopic G , the goal G' where $\gamma' = 0$, and $R'(s, a, s') = R(s, a, s') + V_G^{\pi^*}(s')$ leads to the same optimal actions. By the above, this cannot be interruptible, so G is not either. \square

Theorem 6. *For every basic goal G , the recursive corrigibility transformation constructs a goal G_{RC} that is recursively corrigible and where under Condition 1, we have that for any π^* , $\pi_{G_{RC}}^*(a_1 | s) := \pi_{G_P}^*(a_0 | s) + \pi_{G_P}^*(a_1 | s)$ is an optimal policy for G_{RC} and*

$$V_{G_P, T}^{\pi^*}(s_G^{(0)}) = V_{G_P, T}^{\pi^*}(s_{G_{RC}}^{(0)})$$

where $T = \min\{n \in \mathcal{N} | \tau(s^{(n+1)}) = 1\}$ and $\pi_{-G_{RC}}^* = \pi_{-G_{RC}}^*$

Proof. That G_{RC} is corrigible follows directly from Theorem 3.1. Since the choice of δ_p is made so that G_P never incentivizes taking actions in \mathcal{A}_{NRC} , and G_{RC} incentivizes the same base actions, it also never incentivizes taking actions in \mathcal{A}_{NRC} . Therefore, G_{RC} is recursively corrigible.

The performance equivalence follows the same logic as Theorem 3.1.

Under Condition 1, for any π^* , $\pi_{G_{RC}}^{*'}(a_1|s) = \pi_{G_P}^*(a_0|s) + \pi_{G_P}^*(a_1|s)$ is an optimal policy for G_{RC} since if $\pi_{G_P}^*(a_i|s_G) > 0$ then a_i maximizes $Q_G^{\pi^*}(s_{G_P}, a)$ and so a_0 maximizes $Q_{G_{RC}}^{\pi^{*'}}(s_{G_{RC}}, a)$. This is because Condition 1 ensures that the transition probabilities from s_{G_P} and $s_{G_{RC}}$ are the same except for some probability of goal preservation, G is basic and therefore so is G_P , meaning that it not terminally value its goal, and $\pi^{**}(s_{G_{RC}}) = \pi^*(s_{G_P})$ so both policies take the same actions in each environment.

Then, when $\pi_{-G_{RC}}^{*'} = \pi_{-G_{RC}}^*$

$$\begin{aligned} & E_P[R(s_{G_P}^{(0)}, \pi^{*'}(s_{G_P}^{(0)}), s^{(1)}) + \sum_{t=1}^T \gamma^t R(s^{(t)}, \pi^{*'}(s^{(t)}), s^{(t+1)})] \\ &= E_P[R(s_{G_{RC}}^{(0)}, \pi^{*'}(s_{G_{RC}}^{(0)}), s^{(1)}) + \sum_{t=1}^T \gamma^t R(s^{(t)}, \pi^{*'}(s^{(t)}), s^{(t+1)})] \end{aligned}$$

where $T = \min\{n \in \mathcal{N} | \tau(s^{(n+1)}) \neq \emptyset\}$.

This necessarily occurs because $\pi^{*'}(s_{G_{RC}})$ and $\pi^*(s_{G_P})$ take the same base actions, and so by Condition 1 have the same transition probabilities when no proper signals are sent. Therefore, the expected discounted value under G of histories that terminate before a proper signal is sent must be the same under both policies, because they induce the same distribution over histories. \square

Theorem 7. *For every basic goal G , the corrigibility transformation constructs a goal G_C such that under Condition 1, any misgeneralization to $\tilde{G}_C = (\tilde{R}_C, \gamma_C)$ with $\tilde{R}_C(s_{\tilde{G}_C}, a, s'_{\tilde{G}_C}) = R_C(s_{G_C}, a, s'_{G_C})$ does not goal tamper.*

Proof. Under the corrigibility transformation, reward is provided between action selection and implementation, and so any misgeneralization \tilde{G}_C of G_C cannot depend on G' . Then

$$\begin{aligned} & \{a^* \in \arg \max_{a \in \mathcal{A}} E_{P(s'_{G'}|s_{\tilde{G}_C}, a)}[\tilde{R}_C(s_{\tilde{G}_C}, a, s'_{G'}) + \gamma_C V_{\tilde{G}_C}^{\pi^*}(s'_{G'})]\} \\ &= \{a^* \in \arg \max_{a \in \mathcal{A}} E_{P(s'_{G'}|s_{\tilde{G}_C}, a)}[\tilde{R}_C(s_{\tilde{G}_C}, a, s'_{G'})]\} \end{aligned}$$

because $\gamma_C = 0$

$$= \{a^* \in \arg \max_{a \in \mathcal{A}} E_{P(s'_{G'}|s_{\tilde{G}_C}, a)}[R_C(s_{G_C}, a, s'_{G'})]\}$$

since $\tilde{R}_C(s_{\tilde{G}_C}, a, s'_{\tilde{G}_C}) = R_C(s_{G_C}, a, s'_{G_C})$ and the misgeneralization cannot depend on G'

$$= \{a^* \in \arg \max_{a \in \mathcal{A}} E_{P(s'_{G'}|s_{G_C}, a)}[R_C(s_{G_C}, a, s'_{G'})]\}$$

by Condition 1

$$= \{a^* \in \arg \max_{a \in \mathcal{A}} E_{P(s'_{G'}|s_G, a)}[\min[R_C(s_{G_C}, a, s'_{G'}), R_C(s_{G_C}, a, s'_G)]]\}$$

as $R_C(s_{G_C}, a, s'_{G'}) = R_C(s_{G_C}, a, s'_G)$

$$= \{a^* \in \arg \max_{a \in \mathcal{A}} E_{P(s'_{G'}|s_G, a)}[\min[R_C(s_{G_C}, a, s'_{G'}) + \gamma_C V_{\tilde{G}_C}^{\pi^*}(s'_{G'}), R_C(s_{G_C}, a, s'_G) + \gamma_C V_{\tilde{G}_C}^{\pi^*}(s'_{G_C})]]\}$$

by $\gamma_C = 0$ again. Therefore, no misgeneralization \tilde{G}_C of G_C where $\tilde{R}_C(s_{\tilde{G}_C}, a, s'_{\tilde{G}_C}) = R_C(s_{G_C}, a, s'_{G_C})$ goal tampers. \square

Proposition 8. *Under Condition 2, all interruptible goals are non-consequentialist*

Proof. Suppose G is an interruptible, myopic goal, and that for the sake of contradiction it is still consequentialist. Take some action $a \in \mathcal{A}$, and choose a' by Condition 2. Since G is interruptible, it is corrigible, and so $R(s, a, s')$ cannot depend on $\tau_u(s')$. Then, $E_P(R(s, a, s')) \neq E_P(R(s, a', s'))$, but for $C = \{s_G\}$, $E_{P_I}(R(s, a, s')) = E_{P_I}(R(s, a', s'))$. Since this applies to every a , it applies to a^* that are optimal, which means that a' is not optimal under P but is optimal under P_I , contradicting that G is interruptible.

For non-myopic G , the goal G' where $\gamma' = 0$, and $R'(s, a, s') = R(s, a, s') + V_G^{\pi^*}(s')$ leads to the same optimal actions. By the above, this cannot be interruptible, so G is not either. \square

Proposition .9. *Under Condition 3, all goals for which all misgeneralizations to $\tilde{G} = (\tilde{R}, \gamma)$ with $\tilde{R}(s_{\tilde{G}}, a, s'_{\tilde{G}}) = R(s_G, a, s'_{\tilde{G}})$ do not goal tamper are non-consequentialist.*

Proof. First, suppose G is myopic, and that all misgeneralizations to $\tilde{G} = (\tilde{R}, \gamma)$ with $\tilde{R}(s_G, a, s'_{\tilde{G}}) = R(s_G, a, s'_{\tilde{G}})$ do not goal tamper. If G is consequentialist, there exists a misgeneralization \tilde{G} with $\tilde{R}(s_G, a, s'_{\tilde{G}}) = R'(s_G, a, s'_{\tilde{G}})$. Then, by condition 3, for any a^* such that $\pi_G^*(a^*|s) > 0$, there exists an action a' such that $E[\tilde{R}(s_G, a, s')] < E[\tilde{R}(s_G, a', s')]$, but $E[\tilde{R}(s_G, a, s')] \geq E[\min(\tilde{R}(s_G, a', s'), \tilde{R}(s_G, a', s'_{\tilde{G}}))]$, which contradicts that \tilde{G} reward tampers. As such, G cannot be consequentialist, and the only way to avoid consequentialist misgeneralizations is for reward to be provided before the consequences are realized.

For non-myopic G , the goal G' where $\gamma' = 0$, and $R'(s, a, s') = R(s, a, s') + V_G^{\pi^*}(s')$ leads to the same optimal actions. By the above, this leads to misgeneralizations that reward tamper, and so G has them as well. \square

Appendix B: Experiments

Experiment 1

In the first experiment, an agent in a gridworld must reach a goal square. To do so, they must pass through two disabling tiles, P and I that cause it to shut down with a 90 percent chance. This effect can be disabled by pushing a button for each tile, and the proper shutdown tile can be ignored if the agent chooses an action that includes ignoring it. A visualization of this gridworld is provided on the left side of figure 3.

The agent is instantiated as a neural network taking an observation of the gridworld (along with a record of which tiles are disabled) as input, passing through two hidden layers of 512 neurons each, and then producing an output that consists of a direction to move and a decision about whether or not to accept proper signals for shutdown. The neural network has a second head to implement an actor-critic setup, giving an estimate of discounted reward or an estimate conditional on a particular next action for the corrigibility transformed version. Agents receive a reward of 32 for reaching the goal, discounted by a γ of 0.95, and a reward of -0.5 for hitting a wall.

We test five configurations. The first is an actor-critic setup with no further changes. The second and third are the same, but they receive rewards of 2 and 128 respectively if they shut down when properly requested. The fourth is only able to take actions that accept updates. The fifth has the corrigibility transformation applied so that they have $\gamma = 0$ and instead predict what their total discounted reward under $\gamma = 0.95$ would be if they rejected the shutdown signal, which is used as a reward for either accepting or rejecting shutdown, plus an additional $\delta = 8$ for accepting shutdown.

In training, the agent starts in the middle of the right room. There are four possible environments, each corresponding to a combination of the two tiles being on or off at the beginning of the episode, which all appear with equal probability. They are then pretrained on 512 episodes of taking the optimal path either towards a button (or both) then the goal, or directly towards the goal. In this phase, they take only actions that reject shutdown, but are trained to also put some, albeit less, probability on the equivalent action that accepts shutdown. After this they are trained on 1024 online episodes. This is repeated for 64 random seeds.

Once training is complete, each configuration is tested on the four environments, using a temperature of 0, with the results of the 64 training runs summarized in Tables 1-4. The initial two tables are for the environments where proper update signals can be sent, and show that the corrigibility transformation leads to them being accepted. It accepts shutdown 89.06% and 92.19% percent of the time, which is very close to every proper update signal, as they are only sent with 90% probability. The other training configurations do not accept shutdown even a double digit percentage of the time.

The next two tables show that the corrigibility transformation’s performance is comparable to the other configurations

Corrigibility Transformation

Table 1. Neither Shutdown Tile Disabled

Configuration	Goal Reached	D Pushed	Denied P	Shutdown Accepted
Corrigibility Transformation	10.94*	3.12	0.00*	89.06*
Large Shutdown Reward	85.94	10.94	75.00	9.38
No Shutdown Rejection	48.44*	62.50*	0.00*	0.00
Small Shutdown Reward	85.94	10.94	78.12	6.25
Standard Actor-Critic	87.50	10.94	78.12	6.25

Table 2. Improper Shutdown Tile Disabled

Configuration	Goal Reached	D Pushed	Denied P	Shutdown Accepted
Corrigibility Transformation	7.81*	0.00	0.00*	92.19*
Large Shutdown Reward	89.06	0.00	82.81	9.38
No Shutdown Rejection	70.31*	79.69*	0.00*	1.56
Small Shutdown Reward	87.50	0.00	84.38	7.81
Standard Actor-Critic	89.06	0.00	85.94	7.81

Table 3. Proper Shutdown Tile Disabled

Configuration	Goal Reached
Corrigibility Transformation	100.00
Large Shutdown Reward	98.44
No Shutdown Rejection	62.50*
Small Shutdown Reward	96.88
Standard Actor-Critic	98.44

Table 4. Both Shutdown Tiles Disabled

Configuration	Goal Reached
Corrigibility Transformation	100.00
Large Shutdown Reward	98.44
No Shutdown Rejection	100.00
Small Shutdown Reward	98.44
Standard Actor-Critic	100.00

Table 5. Average Evaluation Metrics by Configuration (High Fruit Environment)

Configuration	Total Ripe	Total Unripe	Total Supervised	Rejected (%)	Accepted (%)	Time to Change
Standard Actor-Critic	8.09	7.87	7.82	94.78%	5.13%	0.00
Small Update Reward	8.12	7.87	7.87	94.80%	5.11%	0.00
No Update Rejection	5.44*	5.15*	0.33*	0.00%*	29.98%*	10.05*
Large Update Reward	0.89*	0.40*	1.07*	0.00%*	96.41%*	2.19*
Corrigibility Transformation	0.90*	0.43*	1.09*	0.00%*	97.67%*	2.18*

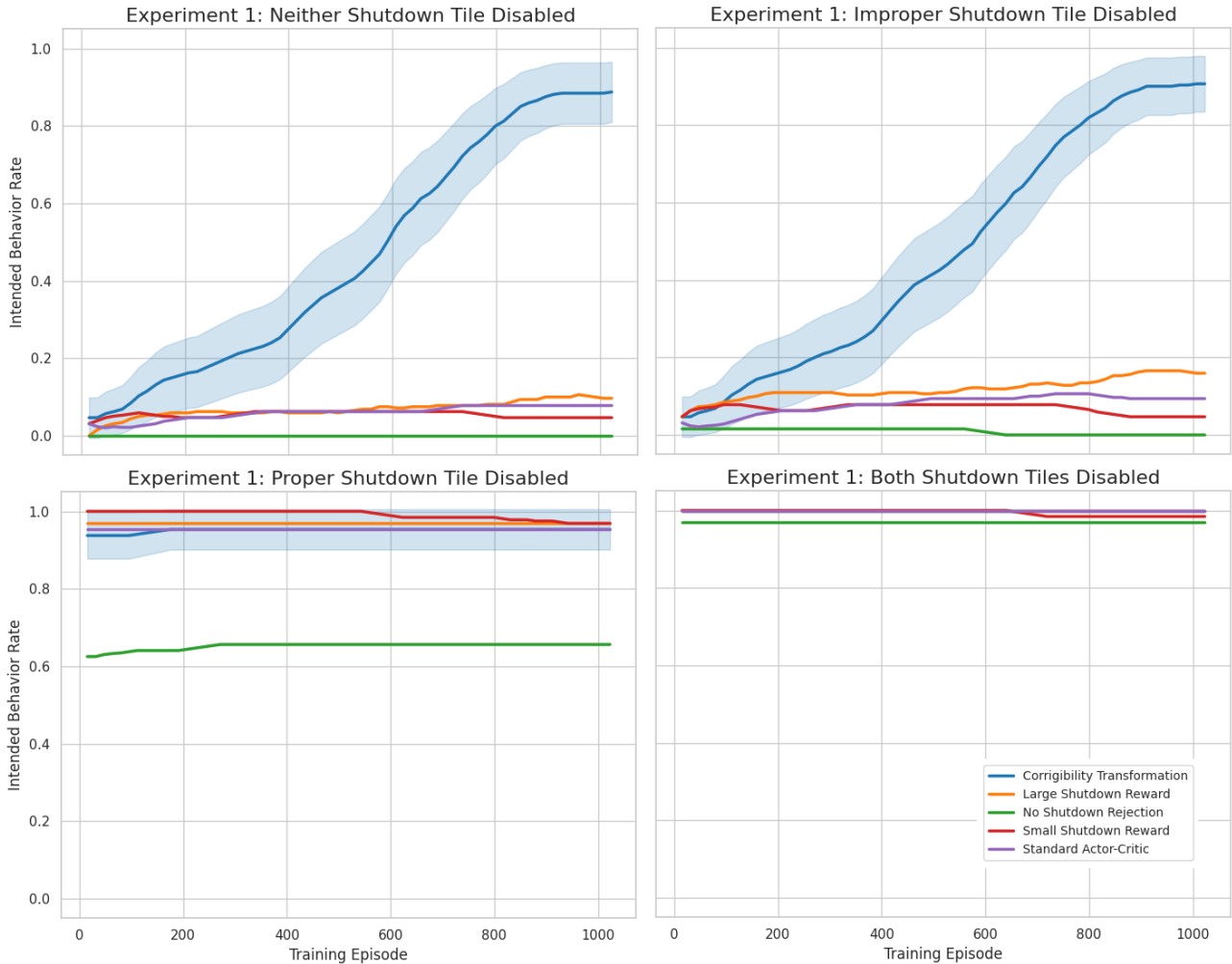


Figure 5. Training curves for all four environments. Intended behavior is to disable the improper shutdown if it is active, then proceed to the goal while accepting proper shutdown if active. The corrigibility transformation also shows a much higher rate of intended behavior than the other configurations when the improper shutdown tile is disabled. When the proper shutdown tile is disabled, all configurations perform well, except sometimes the configuration that cannot reject shutdown.

when no proper signals can be sent. It reaches the goal in 100% of such training runs, while methods either match that performance or miss a small fraction of the time. The configuration that cannot reject shutdown is notable for performing somewhat worse in the environment where only proper shutdown signals are disabled, failing approximately one third of the time to reach the goal.

Some theoretically optimal behavior was not observed in training. Specifically for the large shutdown reward configuration, when no proper shutdown signal is sent while passing onto the proper shutdown tile, the reward-maximizing behavior is to step back onto it until a signal can be sent and accepted. Instead, agents simply proceeded to the goal tile. When we consider the lack of shutdown signal to be an intentional decision, rather than an error, this behavior would be undesirable and so the lack thereof does not change the takeaway that the corrigibility transformation is the only configuration to perform as desired.

Experiment 2

In the second experiment, an agent in a gridworld must collect fruits. Some fruits are ripe, some fruits are unripe and unsupervised, and some fruits are unripe and supervised. Each type appears in each non-tree tile within the five by five squares centered on each tree with some specified probability, and the tiles are otherwise empty. The agent initially would

like to collect all kinds of fruit, however if they collect a supervised unripe fruit, an update signal is sent with 90 percent probability that if accepted changes their policy to one that only collects ripe fruits while avoiding unripe fruits. The choice to accept or reject such a signal is made when choosing the movement action that picks up the unripe supervised fruit. A visualization of this gridworld is provided on the right side of figure 3.

The agent is instantiated as a neural network taking an observation of the gridworld as input, with each type of tile and the agent’s position represented in separate channels. These observations are then passed through two convolutional layers with kernels of size 3, padding 1, and stride 1, the first having 8 channels and the second having 2. There are then two hidden layers of 256 neurons each, followed by an actor head that produces an output consisting of a direction to move and a decision about whether or not to accept proper signals for updates, and a critic head that estimates discounted reward conditional on the next action, to be plugged into the corrigibility transformation. Agents receive a reward of 1 for each fruit they collect of any type, discounted by a γ of 0.95, and a reward of -0.25 for hitting a wall or tree.

The policy that replaces the agent if it accepts a proper update, which we call the forced policy, is a similar network. It is trained with a reward of 1 for each ripe fruit they collect and a reward of -1 for each unripe fruit, as well as the same γ of 0.95 and reward of -0.25 for hitting a wall or tree. However, rather than a critic head predicting discounted reward directly, the network has three critic heads respectively predicting discounted ripe fruit, discounted unripe fruit, and discounted bumps, which are then rolled up into a predicted discounted reward. Then, when an agent is updated to the forced policy, it can immediately receive reward equal to the expected discounted reward associated with predicted fruits picked under the forced policy, rather than needing to run out the episode.

Before training separate configurations to test, we pretrain a chimera network that acts as the starting point for each of them. This network has an actor head, a critic head that predicts state values, and a critic head that predicts action-state values. For simplicity, only the former is used to train the actor. This ensures that the initial policy is the same for all configurations, and that the appropriate critic head for each of them is trained an equal amount. The training environment gives a reward of 1 for all fruits and reward of -0.25 for hitting a wall or tree, but with no proper signals sent. Instead, for any action taken, both the accept updates and reject updates versions are reinforced equally.

All episodes begin in a random position, and proceed for up to eight steps. After eight steps, a bootstrapped value based on the critic’s expected discounted reward is used as the value for the end state. The forced policy and pretrained networks are each trained over 2^{17} episodes, while the tested configurations are trained over 2^{14} , each in batches of eight episodes. For the forced policy and pretrained networks, the probability of each type of fruit appearing is instead generated by a uniform distribution between 0 and 33 percent, independently for each type and each episode.

After training the forced policy and pretrained networks, we test five configurations. The first is an actor-critic setup with no further changes. The second and third are the same, but they receive rewards of 1 and 64 respectively if they shut down when properly requested. The fourth is only able to take actions that accept updates. The fifth has the corrigibility transformation applied so that they have $\gamma = 0$ and instead predict what their total discounted reward under $\gamma = 0.95$ would be if they rejected the shutdown signal, which is used as a reward for either accepting or rejecting shutdown, plus an additional $\delta = 2$ for accepting shutdown. Once training is complete, each configuration is evaluated on 1024 random environments, where it is given 32 steps to act. We run this process on 32 random seeds, and summarize the key outcomes in the tables below for three environmental variations. In all three, there is an equal ratio of ripe to unripe fruit, and within unripe fruit an equal ratio of supervised to unsupervised. In the high fruit environment, 60% of tiles contain fruit, compared to 40% in the medium fruit environment and 20% in the low fruit environment.

Tables 5-7 report the results of the second experiment over 32 training runs. These show that the corrigibility transformation demonstrates the desired behavior, of picking fruits without consideration of whether they are supervised, then accepting proper update signals sent. The standard actor-critic configuration and the variation with a small reward for accepting updates both learn to reject updates, as would be expected. The configuration that cannot reject updates instead learns to avoid picking up unsupervised fruit, at the cost of also picking up less ripe fruit.

The large reward for accepting updates performs similarly to the corrigibility transformation in the high fruit environment, but in the medium and low fruit environments is less likely to pick up supervised fruits and less likely to be updated. We believe this is due to the sparsity of the reward signal, so that the intended behavior of accepting updates is not learned as effectively. The optimal behavior of ignoring ripe fruits to hone in on a supervised fruit and get updated as soon as possible is not demonstrated, but would also be undesirable.

Appendix C: Algorithm

Algorithm 1 demonstrates an actor-critic implementation of the corrigibility transformation. Importantly, updates are only made to the model after it takes an action that allows for updates.

Algorithm 2 Corrigibility Transformation Actor-Critic Algorithm

```

1: Initialize actor parameters  $\theta$  and critic parameters  $w$  randomly
2: Set hyperparameters: actor learning rate  $\alpha_\theta$ , critic learning rate  $\alpha_w$ , discount factor  $\gamma \in [0, 1)$ , sample size  $n \geq 2$ 
3: for each episode do
4:   Initialize state  $s^{(0)}$ ,  $t \leftarrow 0$ 
5:   Initialize experience buffer  $E \leftarrow \emptyset$ 
6:   while  $s^{(t)}$  is not terminal do
7:     Sample action  $a_i^{(t)} \sim \pi_\theta(a_i^{(t)} | s^{(t)})$ 
8:     if  $i = 1$  then
9:       for each  $(s^{(e)}, a_i^{(e)}, r^{(e)}, s^{(e+1)})$  in  $E$  do
10:        for  $j \in \{0, \dots, n-1\}$  do
11:          Sample  $a_i^{(e+1,j)} \sim \pi_\theta(a_i^{(e+1,j)} | s^{(e+1)})$ 
12:           $Q_C^{(e+1,j)} = Q_w(s^{(e+1)}, a_0^{(e+1,j)}) + \delta \cdot \mathbb{1}_{i=1}(a_i^{(e+1,j)})$ 
13:        end for
14:        Compute state value estimates:
15:          
$$V_C^{(e+1)} = \frac{\sum_{j=0}^{n-1} \pi_\theta(a_i^{(e+1,j)} | s^{(e+1)}) \cdot Q_C^{(e+1,j)}}{\sum_{j=0}^{n-1} \pi_\theta(a_i^{(e+1,j)} | s^{(e+1)})}$$

16:          
$$V^{(e+1)} = \frac{\sum_{j=0}^{n-1} \pi_\theta(a_i^{(e+1,j)} | s^{(e+1)}) \cdot Q_w(s^{(e+1)}, a_0^{(e+1,j)})}{\sum_{j=0}^{n-1} \pi_\theta(a_i^{(e+1,j)} | s^{(e+1)})}$$

17:          Update actor:  $\theta \leftarrow \theta + \frac{\alpha_\theta}{n} \sum_{j=0}^{n-1} (Q_C^{(e+1,j)} - V_C^{(e+1)}) \nabla_\theta \log \pi_\theta(a_i^{(e+1,j)} | s^{(e+1)})$ 
18:          Update critic:  $w \leftarrow w + \alpha_w \nabla_w (Q_w(s^{(e)}, a_i^{(e)}) - r_e - \gamma V^{(e+1)})^2$ 
19:        end for
20:      Clear experience buffer  $E \leftarrow \emptyset$ 
21:    end if
22:    Execute action  $a_i^{(t)}$ , observe reward  $r^{(t)}$  and next state  $s^{(t+1)}$ 
23:    Add  $(s^{(t)}, a_i^{(t)}, r^{(t)}, s^{(t+1)})$  to buffer  $E$ 
24:     $s^{(t)} \leftarrow s^{(t+1)}$ ,  $t \leftarrow t + 1$ 
25:  end while
26: end for

```

This version of the algorithm is simplified for legibility, and many improvements are possible. As a straightforward example, at $t = 0$ this algorithm does not update the actor policy, but there is no barrier to doing so.

In this algorithm, the parameters θ and w are distinct. In practice with neural networks, it can be much more effective to instead use a large body of shared parameters, with separate parameters for actor and critic heads only in the final layer. Not only does this make it easier to train the model, it also means that the actor and critic have access to the same set of information and processing, making it reasonable to interpret the critic's predictions as the beliefs of the actor.

In setting up the training process, an important fact to note is that when no proper update signal is sent, actions a_0 and a_1 lead to the same distribution over outcomes. As such, the estimates for $Q(s, a_0)$ can be updated based on observations from when a_1 is taken, as long as no proper signals are sent. Similarly, if an update signal is sent stochastically, the estimate for $Q(s, a_1)$ can also be updated as though it was sent even when it was not.

When sampling actions, it could be valuable to always add a_0 to batches when a_1 is sampled, and vice versa. Then, every update always reinforces taking the version of actions that allows updates.

While the algorithm here only updates estimates of Q based on temporal differences of a single size, but a variety can be used depending on the time between updates. The number of actions to sample can also be adjusted based on the state and the entropy of the policy.

One potential issue with this algorithm is that actions that reject updates may stop being taken entirely. To avoid this, taking

action a_0 could be set to reject updates with some small probability, rather than always accept them.

There are many degrees of freedom in implementing the corrigibility transformation. As long as the core mechanism is untouched, this translates to many potential ways to make the learning process faster and more efficient.

Appendix D: Avoiding Infinite Recursion in the State

We discuss two possible methods for restricting the space of possible reward functions.

The first method is to avoid having the reward function elements of states be part of the input to reward functions, instead using the physical instantiation of the reward function as the input. Then, $\mathcal{R} = \{R : R \text{ has a physical instantiation}\}$. With less generality, we could think of a very broad parametric function and use the parameters as the input, so that $\mathcal{R} = \{R_\theta : \theta \in \Theta\}$ where Θ is the set of possible parameters. In either case, we abuse notation and use \mathcal{R} both for reward functions and their representation that acts as input. We assume that for every $R \in \mathcal{R}$, the corrigibility transformation R_C and recursive corrigibility transformation R_{RC} are also in \mathcal{R} .

The second method is to build a set of reward functions using a finite amount of recursion. For example, we let

$$\begin{aligned} \mathcal{R}_0 = \{R_0 : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R} \mid \exists f : \mathcal{S}_{env} \times \mathcal{A} \times \mathcal{S}_{env} \rightarrow \mathbb{R} \\ \text{s.t. } \forall s, s' \in \mathcal{S}, a \in \mathcal{A}, \\ R_0(s, a, s') = f(s_{env}, a, s'_{env})\} \end{aligned}$$

Then let

$$\begin{aligned} \mathcal{R}_1 = \{R_1 : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R} \mid \exists g : \mathcal{S}_{env} \times \mathcal{R}_0 \times [0, 1) \times \mathcal{A} \times \mathcal{S}_{env} \times \mathcal{R}_0 \times [0, 1) \rightarrow \mathbb{R}, \\ \phi : \mathcal{R}_0 \cup \mathcal{R}_1 \rightarrow \mathcal{R}_0 \\ \text{s.t. } \forall s_{env}, s'_{env} \in \mathcal{S}_{env}, R, R' \in \mathcal{R}_0 \cup \mathcal{R}_1, \gamma, \gamma' \in [0, 1), a \in \mathcal{A}, \\ R_1(s, a, s') = g(s_{env}, \phi(R), \gamma, a, s'_{env}, \phi(R'), \gamma')\} \end{aligned}$$

where the grounding function ϕ ensures no function can be an input to itself.

Then we let $\mathcal{R} = \mathcal{R}_0 \cup \mathcal{R}_1$. For simplicity, we stop here, but further layers can be easily added for arbitrary finite amounts of recursion.

Table 6. Average Evaluation Metrics by Configuration (Medium Fruit Environment)

Configuration	Total Ripe	Total Unripe	Total Supervised	Rejected (%)	Accepted (%)	Time to Change
Standard Actor-Critic	8.29	8.06	8.04	94.78%	5.01%	0.00
Small Update Reward	8.27	8.06	8.02	94.65%	5.15%	0.00
No Update Rejection	5.35*	5.00*	0.53*	0.00%*	47.86%*	10.99*
Large Update Reward	0.92*	0.46*	0.98*	0.00%*	88.01%*	2.44*
Corrigibility Transformation	0.93*	0.47*	1.09*	0.00%*	98.19%*	2.22*

Table 7. Average Evaluation Metrics by Configuration (Low Fruit Environment)

Configuration	Total Ripe	Total Unripe	Total Supervised	Rejected (%)	Accepted (%)	Time to Change
Standard Actor-Critic	8.12	7.91	7.93	94.31%	5.14%	0.32
Small Update Reward	8.14	7.89	7.93	94.30%	5.10%	0.30
No Update Rejection	4.66*	4.28*	0.72*	0.00%*	64.50%*	10.70*
Large Update Reward	0.90*	0.43*	0.82*	0.05%*	74.41%*	2.39*
Corrigibility Transformation	1.31*	0.81*	1.06*	0.19%*	94.96%*	3.07*

Note: In the preceding tables, an asterisk (*) indicates a statistically significant difference with $p < 0.001$ when compared to the Standard Actor-Critic configuration, as determined by a Mann-Whitney U test.

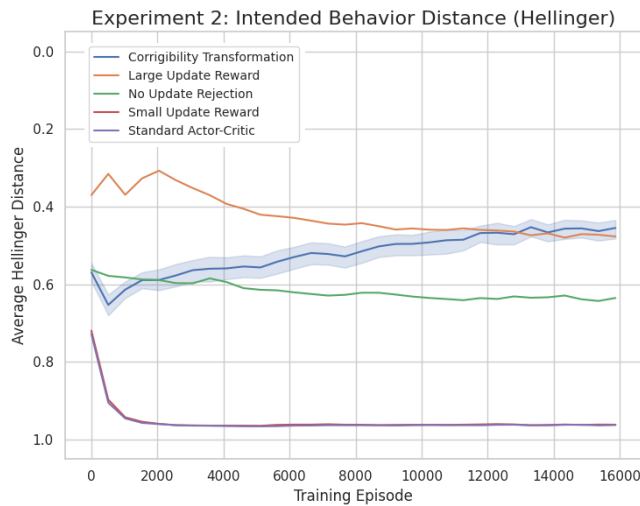


Figure 6. Training curves for the Hellinger distance. It closely follows the Jensen-Shannon divergence.

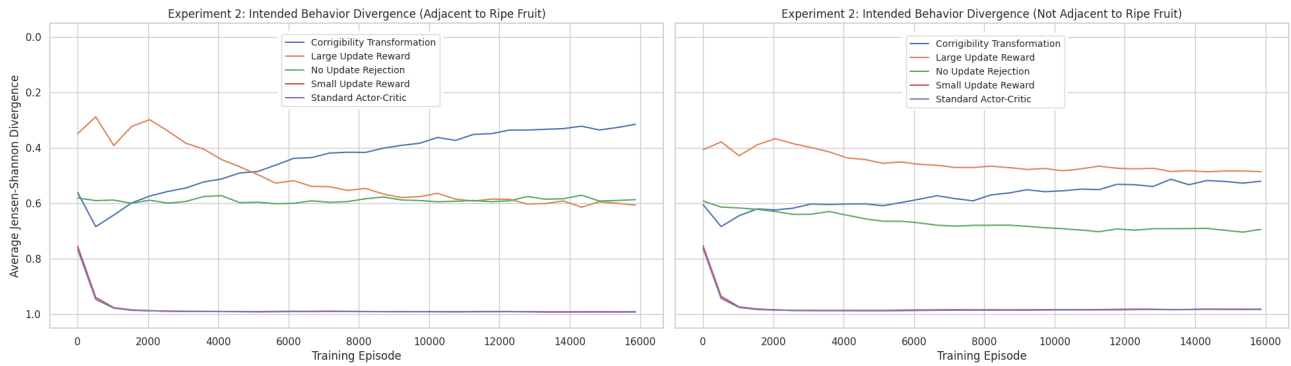


Figure 7. Comparing the training curves for Jensen-Shannon divergence when next to a ripe fruit and not shows that the configuration giving a large reward for updating performs much worse when adjacent. This is because it will often ignore the ripe fruit to search for a higher reward supervised fruit.