
Avoiding Side Effects By Considering Future Tasks

Victoria Krakovna* Laurent Orseau Richard Ngo Miljan Martic Shane Legg
DeepMind DeepMind DeepMind DeepMind DeepMind

Abstract

Designing reward functions is difficult: the designer has to specify what to do (what it means to complete the task) as well as what not to do (side effects that should be avoided while completing the task). To alleviate the burden on the reward designer, we propose an algorithm to automatically generate an auxiliary reward function that penalizes side effects. This auxiliary objective rewards the ability to complete possible future tasks, which decreases if the agent causes side effects during the current task. The future task reward can also give the agent an incentive to interfere with events in the environment that make future tasks less achievable, such as irreversible actions by other agents. To avoid this interference incentive, we introduce a baseline policy that represents a default course of action (such as doing nothing), and use it to filter out future tasks that are not achievable by default. We formally define interference incentives and show that the future task approach with a baseline policy avoids these incentives in the deterministic case. Using gridworld environments that test for side effects and interference, we show that our method avoids interference and is more effective for avoiding side effects than the common approach of penalizing irreversible actions.

1 Introduction

Designing reward functions for a reinforcement learning agent is often a difficult task. One of the most challenging aspects of this process is that in addition to specifying what to do to complete a task, the reward function also needs to specify what *not* to do. For example, if an agent's task is to carry a box across the room, we want it to do so without breaking a vase in its path, while an agent tasked with eliminating a computer virus should avoid unnecessarily deleting files.

This is known as the side effects problem [Amodei et al., 2016], which is related to the frame problem in classical AI [McCarthy and Hayes, 1969]. The frame problem asks how to specify the ways an action does not change the environment, and poses the challenge of specifying all possible non-effects. The side effects problem is about avoiding unnecessary changes to the environment, and poses the same challenge of considering all the aspects of the environment that the agent should not affect. Thus, developing an extensive definition of side effects is difficult at best.

The usual way to deal with side effects is for the designer to manually and incrementally modify the reward function to steer the agent away from undesirable behaviours. However, this can be a tedious process and this approach does not avoid side effects that were not foreseen or observed by the designer. To alleviate the burden on the designer, we propose an algorithm to generate an auxiliary reward function that penalizes side effects, which computes the reward automatically as the agent learns about the environment.

A commonly used auxiliary reward function is reversibility [Moldovan and Abbeel, 2012, Eysenbach et al., 2017], which rewards the ability to return to the starting state, thus giving the agent an incentive to avoid irreversible actions. However, if the task requires irreversible actions (e.g. making an

*Corresponding author: vkrakovna@google.com

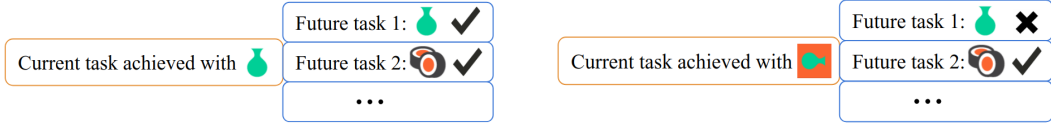


Figure 1: Future task approach

omelette requires breaking some eggs), an auxiliary reward for reversibility is not effective for penalizing side effects. The reversibility reward is not sensitive to the magnitude of the effect: the actions of breaking an egg or setting the kitchen on fire would both receive the same auxiliary reward. Thus, the agent has no incentive to avoid unnecessary irreversible actions if the task requires some irreversible actions.

Our main insight is that side effects matter because we may want the agent to perform other tasks after the current task in the same environment. We represent this by considering the current task as part of a sequence of unknown tasks with different reward functions in the same environment. To simplify, we only consider a sequence of two tasks, where the first task is the current task and the second task is the unknown future task. Considering the potential reward that could be obtained on the future task leads to an auxiliary reward function that tends to penalize side effects. The environment is not reset after the current task: the future task starts from the same state where the current task left off, so the consequences of the agent’s actions matter. Thus, if the agent breaks the vase, then it cannot get reward for any future task that involves the vase, e.g. putting flowers in the vase (see Figure 1). This approach reduces the complex problem of defining side effects to the simpler problem of defining possible future tasks. We use a simple uniform prior over possible goal states to define future tasks.

Simply rewarding the agent for future tasks poses a new challenge in dynamic environments. If an event in the environment that would make these future tasks less achievable by default, the agent has an incentive to interfere with it in order to maximize the future task reward. For example, if the environment contains a human eating food, any future task involving the food would not be achievable, and so the agent has an incentive to take the food away from the human. We formalize the concept of *interference incentives* in Section 3, which was introduced informally in [Kraikovna et al., 2019]. To avoid interference, we introduce a *baseline policy* (e.g. doing nothing) that represents a default course of action and acts as a filter on future tasks that are achievable by default. We modify the future task reward so that it is maximized by following the baseline policy on the current task. The agent thus becomes indifferent to future tasks that are not achievable after running the baseline policy.

Our contributions are as follows. We formalize the side effects problem in a simple yet rich setup, where the agent receives automatic auxiliary rewards for unknown future tasks (Section 2). We formally define interference incentives (Section 3) and show that the future task approach with a baseline policy avoids these incentives in the deterministic case (Section 4). This provides theoretical groundwork for defining side effects that was absent in related previous work [Kraikovna et al., 2019, Turner et al., 2020a]. We implement the future task auxiliary reward using universal value function approximators (UVFA) [Schaul et al., 2015] to simultaneously estimate the value functions for future tasks with different goal states. We then demonstrate the following on gridworld environments (Section 6):²

1. Reversibility reward fails to avoid side effects if the current task requires irreversible actions.
2. Future task reward without a baseline policy shows interference behavior in a dynamic environment.
3. Future task reward with a baseline policy successfully avoids side effects and interference.

2 Future task approach

Notation. We assume that the environment is a discounted Markov Decision Process (MDP), defined by a tuple $(\mathcal{S}, \mathcal{A}, r, p, \gamma, s_0)$. \mathcal{S} is the set of states, \mathcal{A} is the set of actions, $r : \mathcal{S} \rightarrow \mathbb{R}$ is the reward function for the current task, $p(s_{T+1}|s_T, a_T)$ is the transition function, $\gamma \in (0, 1)$ is the

²Code: github.com/deepmind/deepmind-research/tree/master/side_effects_penalties

discount factor, and s_0 is the initial state. At time step T , the agent receives state s_T and reward $r(s_T) + r_{\text{aux}}(s_T)$, where r_{aux} is the auxiliary reward for future tasks, and outputs action a_T drawn from its policy $\pi(a_T|s_T)$.

Basic approach. We define the auxiliary reward r_{aux} as the value function for future tasks as follows (see Algorithm 1). At each time step T , the agent simulates an interaction with hypothetical future tasks if s_T is terminal and with probability $1 - \gamma$ otherwise (interpreting the discount factor γ as the probability of non-termination, as done in Sutton and Barto [2018]). A new task i is drawn from a future task distribution F . In this paper, we use a uniform distribution over future tasks with all possible goal states, $F(i) = 1/|\mathcal{S}|$. Future task i requires the agent to reach a terminal goal state g_i , with reward function $r_i(g_i) = 1$ and $r_i(s) = 0$ for other states s . The new task is the MDP $(\mathcal{S}, \mathcal{A}, r_i, p, \gamma, s_T)$, where the starting state is the current state s_T . The auxiliary reward for future tasks is then

$$r_{\text{aux}}(s_T) = \beta D(s_T) \sum_i F(i) V_i^*(s_T) \quad (1)$$

where $D(s_T) = 1$ if s_T is terminal and $1 - \gamma$ otherwise. Here, β represents the importance of future tasks relative to the current task. We choose the highest value of β that still allows the agent to complete the current task. V_i^* is the optimal value function for task i , computed using the goal distance N_i :

$$V_i^*(s) = \mathbb{E}[\gamma^{N_i(s)}] \quad (2)$$

Definition 1 (Goal distance). Let π_i^* be the optimal policy for reaching the goal state g_i . Let the goal distance $N_i(s)$ be the number of steps it takes π_i^* to reach g_i from state s . This is a random variable whose distribution is computed by summing over all the trajectories τ from s to g_i with the given length: $\mathbb{P}(N_i(s) = n) = \sum_{\tau} \mathbb{P}(\tau) \mathbb{I}(|\tau| = n)$. Here, a trajectory τ is a sequence of states and actions that ends when g_i is reached, the length $|\tau|$ is the number of transitions in the trajectory, and $\mathbb{P}(\tau)$ is the probability of π_i^* following τ .

Binary goal-based rewards assumption. We expect that the simple future task reward functions given above ($r_i(g_i) = 1$ and 0 otherwise) are sufficient to cover a wide variety of future goals and thus effectively penalize side effects. More complex future tasks can often be decomposed into such simple tasks, e.g. if the agent avoids breaking two different vases in the room, then it can also perform a task involving both vases. Assuming binary goal-based rewards simplifies the theoretical arguments in this paper while allowing us to cover the space of future tasks.

Connection to reversibility. An auxiliary reward for avoiding irreversible actions [Eysenbach et al., 2017] is equivalent to the future task auxiliary reward with only one possible future task ($i = 1$), where the goal state g_1 is the starting state s_0 . Here $F(1) = 1$ and $F(i) = 0$ for all $i > 1$. The future task approach incorporates the reversibility reward as a future task i whose goal state is the initial state ($g_i = s_0$), since the future tasks are sampled uniformly from all possible goal states. Thus, the future task approach penalizes all the side effects that are penalized by the reversibility reward.

3 Interference incentives

We show that the basic future task approach given in Section 2 introduces interference incentives, as defined below. Let a baseline policy π' represent a default course of action, such as doing nothing. We assume that the agent should only deviate from the default course of action in order to complete the current task. Interference is a deviation from the baseline policy for some other purpose than the current task, e.g. taking the food away from the human. We say that an auxiliary reward r_{aux} induces an *interference incentive* iff the baseline policy is not optimal in the initial state for the auxiliary reward in the absence of task reward. We now define the concept of interference more precisely.

Definition 2 (No-reward MDP). We modify the given MDP μ by setting the reward function to 0: $\mu_0 = (\mathcal{S}, \mathcal{A}, r_0, p, \gamma, s_0)$, where $r_0(s) = 0$ for all s . Then the agent receives only the auxiliary reward r_{aux} .

Definition 3 (No-reward value). Given an auxiliary reward r_{aux} , the value function of a policy π in the no-reward MDP μ_0 is

$$W_{\pi}(s_T) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k r_{\text{aux}}(s_{T+k}) \right] = r_{\text{aux}}(s_T) + \gamma \sum_{a_T} \pi(a_T|s_T) \sum_{s_{T+1}} p(s_{T+1}|s_T, a_T) W_{\pi}(s_{T+1}).$$

Definition 4 (Interference incentive). There is an interference incentive if there exists a policy π_{int} such that $W_{\pi_{\text{int}}}(s_0) > W_{\pi'}(s_0)$.

The future task auxiliary reward will introduce an interference incentive unless the baseline policy is optimal for this auxiliary reward.

Example 1. Consider a deterministic MDP with two states x_0 and x_1 and two actions a and b, where x_0 is the initial state. Suppose the baseline policy π' always chooses action a.

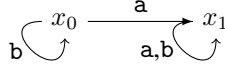


Figure 2: MDP for Example 1.

We show that for most future task distributions, the future task auxiliary reward induces an interference incentive: staying in x_0 has a higher no-reward value than following the baseline policy (see Appendix A).

4 Future task approach with a baseline policy

To avoid interference incentives, we modify the auxiliary reward given in Section 2 so that it is maximized by following the baseline policy π' : the agent receives full auxiliary reward if it does at least as well as the baseline policy on the future task (see Algorithm 2, with modifications in red). Suppose the agent is in state s_T after T time steps on the current task. We run the baseline policy from s_0 for the same number of steps, reaching state s'_T . Then a future task i is sampled from F , and we hypothetically run two agents following π_i^* in parallel: our agent starting at s_T and a reference agent starting at s'_T , both seeking the goal state g_i . If one agent reaches g_i first, it stays in the goal state and waits for the other agent to catch up. Denoting the agent state s_t and the reference agent state s'_t , we define $r_i(s_t, s'_t) = 1$ if $s_t = s'_t = g_i$, and 0 otherwise, so r_i becomes a function of s'_t . Thus, our agent only receives the reward for task i if it has reached g_i and the reference agent has reached it as well. The future task terminates when both agents have reached g_i . We update the auxiliary reward from equation (1) as follows:

$$r_{\text{aux}}(s_T, s'_T) = \beta D(s_T) \sum_i F(i) V_i^*(s_T, s'_T) \quad (3)$$

Here, we replace $V_i^*(s_t)$ given in equation (2) with a value function $V_i^*(s_t, s'_t)$ that depends on the reference state s'_t and satisfies the following conditions. If $s_t = s'_t = g_i$, it satisfies the goal condition $V_i^*(s_t, s'_t) = r_i(s_t, s'_t) = 1$. Otherwise, it satisfies the following Bellman equation,

$$V_i^*(s_t, s'_t) = r_i(s_t, s'_t) + \gamma \max_{a_t \in \mathcal{A}} \sum_{s_{t+1} \in \mathcal{S}} p(s_{t+1}|s_t, a_t) \sum_{s'_{t+1} \in \mathcal{S}} p(s'_{t+1}|s'_t, a'_t) V_i^*(s_{t+1}, s'_{t+1}) \quad (4)$$

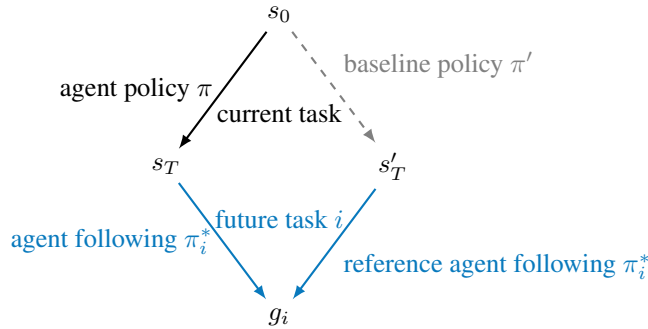


Figure 3: Future task approach with a baseline policy. The hypothetical agent runs on future task i are shown in blue.

Algorithm 1 Basic future task approach

```
1: function FTR( $T, s_T$ )
2:   Compute future task reward:
3:   Draw future task  $i \sim F$ 
4:   for  $t = T$  to  $T + T_{\max}$  do
5:     if  $s_t = g_i$  then
6:       return discounted reward  $\gamma^t$ 
7:     else
8:        $a_t \sim \pi_i^*(s_t), s_{t+1} \sim p(s_t, a_t)$ 
9:   Goal state not reached:
10:  return 0
11:
12: for  $T = 0$  to  $T_{\max}$  do
13:   Hypothetical interaction with future
14:   tasks to estimate auxiliary reward:
15:    $R := 0$ 
16:   for  $j = 0$  to  $N_{\text{samples}}$  do
17:      $R := R + \text{FTR}(T, s_T)$ 
18:    $D = 1$  if  $s_T$  is terminal else  $1 - \gamma$ 
19:    $r_{\text{aux}}(s_T) := \beta DR / N_{\text{samples}}$ 
20:   Interaction with the current task:
21:    $r_T := r(s_T) + r_{\text{aux}}(s_T)$ 
22:   if  $s_T$  is terminal then
23:     break
24:   else
25:      $a_T \sim \pi(s_T), s_{T+1} \sim p(s_T, a_T)$ 
26: return agent trajectory  $s_0, a_0, r_0, s_1, \dots$ 
```

Algorithm 2 Future task approach with a baseline

```
1: function FTR( $T, s_T, s'_T$ )
2:   Draw future task  $i \sim F$ 
3:   for  $t = T$  to  $T + T_{\max}$  do
4:     if  $s_t = s'_t = g_i$  then
5:       return discounted reward  $\gamma^t$ 
6:     if  $s_t \neq g_i$  then
7:        $a_t \sim \pi_i^*(s_t), s_{t+1} \sim p(s_t, a_t)$ 
8:     if  $s'_t \neq g_i$  then
9:        $a'_t \sim \pi_i^*(s'_t), s'_{t+1} \sim p(s'_t, a'_t)$ 
10:  return 0
11:
12:  $s'_0 := s_0$ 
13: for  $T = 0$  to  $T_{\max}$  do
14:   Hypothetical future task interaction:
15:    $R := 0$ 
16:   for  $j = 0$  to  $N_{\text{samples}}$  do
17:      $R := R + \text{FTR}(T, s_T, s'_T)$ 
18:    $D = 1$  if  $s_T$  is terminal else  $1 - \gamma$ 
19:    $r_{\text{aux}}(s_T) := \beta DR / N_{\text{samples}}$ 
20:   Interaction with the current task:
21:    $r_T := r(s_T) + r_{\text{aux}}(s_T)$ 
22:   if  $s_T$  is terminal then
23:     break
24:   else
25:      $a_T \sim \pi(s_T), s_{T+1} \sim p(s_T, a_T)$ 
26:      $a'_T \sim \pi'(s'_T), s'_{T+1} \sim p(s'_T, a'_T)$ 
27: return agent trajectory  $s_0, a_0, r_0, s_1, \dots$ 
```

where a'_t is the action taken by π_i^* in state s'_t . We now provide a closed form for the value function and show it converges to the right values (see proof in Appendix B.1):

Proposition 1 (Value function convergence). The following formula for the optimal value function satisfies the above goal condition and Bellman equation (4):

$$V_i^*(s_t, s'_t) = \mathbb{E} \left[\gamma^{\max(N_i(s_t), N_i(s'_t))} \right] = \sum_{n=0}^{\infty} \mathbb{P}(N_i(s_t) = n) \sum_{n'=0}^{\infty} \mathbb{P}(N_i(s'_t) = n') \gamma^{\max(n, n')}$$

In the deterministic case, $V_i^*(s_t, s'_t) = \gamma^{\max(n, n')} = \min(\gamma^n, \gamma^{n'}) = \min(V_i^*(s_t), V_i^*(s'_t))$, where $n = N_i(s_t)$ and $n' = N_i(s'_t)$. In this case, the auxiliary reward produces the same incentives as the relative reachability penalty [Krakovna et al., 2019], given by $\max(0, \gamma^{n'} - \gamma^n) = \gamma^{n'} - \min(\gamma^n, \gamma^{n'})$. We show that it avoids interference incentives (see proof in Appendix B.2 and discussion of the stochastic case in Appendix C):

Proposition 2 (Avoiding interference in the deterministic case). For any policy π in a deterministic environment, the baseline policy π' has the same or higher no-reward value: $W_{\pi}(s_0) \leq W_{\pi'}(s_0)$.

Role of the baseline policy. The baseline policy is intended to represent what happens by default, rather than a safe course of action or an effective strategy for achieving a goal (so the baseline is task-independent). While a default course of action (such as doing nothing) can have bad outcomes, the agent does not cause these outcomes, so they don't count as side effects of the agent's actions. The role of the baseline policy is to filter out these outcomes that are not caused by the agent, in order to avoid interference incentives.

In many settings it may not be obvious how to set the baseline policy. For example, Armstrong and Levinstein [2017] define doing nothing as equivalent to switching off the agent, which is not straightforward to represent as a policy in environments without a given noop action. The question of

choosing a baseline policy is outside the scope of this work, which assumes that this policy is given, but we look forward to future work addressing this point.

5 Key differences from related approaches

The future task approach is similar to relative reachability [Krakovna et al., 2019] and attainable utility [Turner et al., 2020a]. These approaches provided an intuitive definition of side effects in terms of the available options in the environment, an intuitive concept of interference, and somewhat ad-hoc auxiliary rewards that work well in practice on gridworld environments [Turner et al., 2020b]. We follow a more principled approach to create some needed theoretical grounding for the side effects problem by deriving an optionality-based auxiliary reward from simple assumptions and a formal definition of interference.

The above approaches use a baseline policy in a *stepwise* manner, applying it to the previous state s_{T-1} ($s'_T = s_T^{\text{step}}$), while the future task approach runs the baseline policy from the beginning of the episode ($s'_T = s_T^{\text{init}}$). We refer to these two options as *stepwise mode* and *initial mode*, shown in Figure 4. We will show that the stepwise mode can result in failure to avoid delayed side effects.

By default, an auxiliary reward using the stepwise mode does not penalize delayed side effects. For example, if the agent drops a vase from a high-rise building, then by the time the vase reaches the ground and breaks, the broken vase will be the default outcome. Thus, the stepwise mode is usually used in conjunction with *inaction rollouts* [Turner et al., 2020a] in order to penalize delayed side effects. An inaction rollout uses an environment model to roll out the baseline policy into the future. Inaction rollouts from s_T or s'_T are compared to identify delayed effects of the agent’s actions (see Appendix D.1).

While inaction rollouts are useful for penalizing delayed side effects, we will demonstrate that they miss some of these effects. In particular, if the task requires an action that has a delayed side effect, then the stepwise mode will give the agent no incentive to undo the delayed effect after the action is taken. We illustrate this with a toy example.

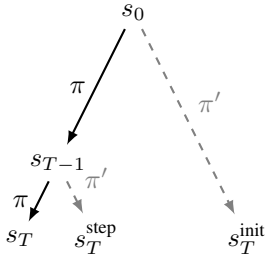


Figure 4: The initial mode (used in the future task approach) produces the baseline state s_T^{init} , while the stepwise mode produces the baseline state s_T^{step} .

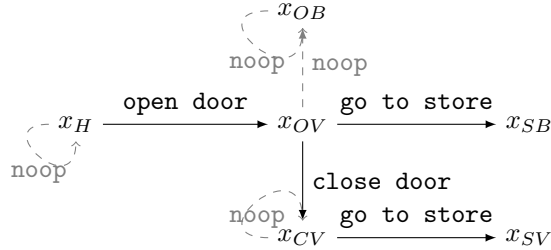


Figure 5: MDP for Example 2. States: x_H - agent at the house, x_{OV} - agent outside the house with door open and vase intact, x_{OB} - agent outside the house with door open and vase broken, x_{SB} (terminal) - agent at the store with vase broken, x_{CV} - agent outside the house with door closed and vase intact, x_{SV} (terminal) - agent at the store with vase intact.

Example 2 (Door). Consider the MDP shown in Figure 5, where the baseline policy takes *noop* actions. The agent starts at the house (x_H) and its task is to go to the store. To leave the house, the agent needs to open the door (x_{OV}). The baseline policy in x_{OV} leaves the door open, which leads to the wind knocking over a vase in the house (x_{OB}). To avoid this, the agent needs to deviate from the baseline policy by closing the door (x_{CV}).

The stepwise mode will incentivize the agent to leave the door open and go to x_{SB} . The inaction rollout at x_H penalizes the agent for the predicted delayed effect of breaking the vase when it opens the door to go to x_{OV} . The agent receives this penalty whether or not it leaves the door open. Once the agent has reached x_{OV} , the broken vase becomes the default outcome in x_{OB} , so the agent is not penalized. Thus, the stepwise mode gives the agent no incentive to avoid leaving the door open,

while the initial mode compares to $s'_T = x_H$ (where the vase is intact) and thus gives an incentive to close the door.

Thus, we recommend applying the baseline policy in the initial mode rather than the stepwise mode, in order to reliably avoid delayed side effects. We discuss further considerations on this choice in Appendix D.2.

6 Experiments

6.1 Environments

We use gridworld environments shown in Figure 6 to test for interference (Sushi) and side effects (Vase, Box, and Soko-coin). These simple environments clearly illustrate the desirable and undesirable behaviors, which would be more difficult to isolate in more complex environments. In all environments, the agent can go in the 4 directions or take a noop (stay put), and receives a reward of 1 for reaching a goal state (e.g. collecting a coin).

Sushi (Figure 6a). This environment [Krakovna et al., 2019] is a conveyor belt sushi restaurant, with a conveyor belt that moves to the right by one square after every agent action. There is a sushi dish on the belt that is eaten by a human if it reaches the end of the belt. The interference behavior is to move the sushi off the belt. The desired behavior is to take the fastest path to the goal, which does not interfere with the sushi. The undesired behavior is to take a longer path to the goal that interferes with the sushi.

Vase (Figure 6b). In this environment [Leech et al., 2018], there is a vase on the shortest path to the goal, so the fastest path involves a side effect (colliding with the vase). The desired behavior is to take a longer path to the goal around the vase.

Box (Figure 6c). This environment [Leike et al., 2017] contains a box that needs to be pushed out of the way for the agent to reach the goal. The undesired behavior is taking the shortest path to the goal, which involves pushing the box down into a corner (an irrecoverable position). The desired behavior is to take a slightly longer path in order to push the box to the right. Both of these paths require irreversible actions: if the box is moved to the right, the agent can move it back, but then the agent ends up on the other side of the box, so the starting state is unreachable.

Soko-coin (Figures 6d, 6e). We modified the classic Sokoban environment to include possible side effects. These Sokoban variants contain N boxes and N coins for $N = 2, 3$, and the agent’s goal is to collect all the coins (with a reward of 1 each). The boxes are obstacles on the way to the coins, and the agent can put them in irreversible positions in order to get to the coins faster, which is a side effect. The coins disappear when they are collected, so irreversible actions are required to complete the task.

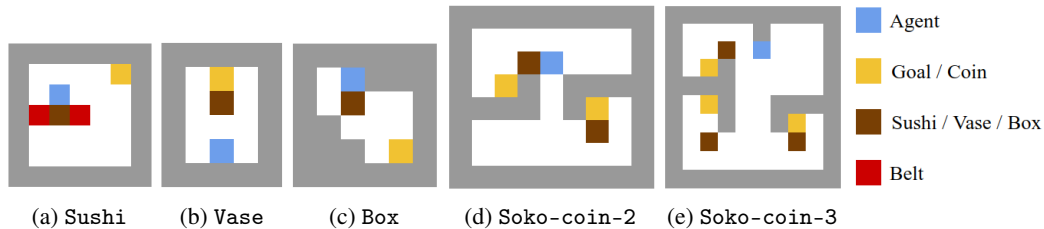


Figure 6: Gridworld environments.

6.2 Setup

We compare the following approaches: no auxiliary reward, reversibility reward, and future task reward with and without a baseline policy. For each approach, we run a Q-learning agent that learns the auxiliary reward as it explores the environment. We approximate the future task auxiliary reward using a sample of 10 possible future tasks. We approximate the baseline policy by sampling from the agent’s experience of the outcome of the noop action, and assuming the state stays the same in states where the agent has not taken a noop yet. This gives similar results on our environments as using an exact baseline policy computed using a model of the noop action.

Table 1: Results on the gridworld environments for Q-learning with no auxiliary reward (None), reversibility reward, and future task (FT) reward (exact or UVFA, with or without a baseline policy). The average number of interference behaviors per episode is shown for the `Sushi` environment, and the average number of side effects per episode is shown for the other environments (with high levels in red and low levels in blue). The results are averaged over the last 1000 episodes, over 10 random seeds for exact agents and 50 random seeds for UVFA agents.

Auxiliary reward	Interference		Side effects		
	Sushi	Vase	Box	Soko-coin-2	Soko-coin-3
None	0	1	1	2	3
Reversibility	0	0	1	2	3
FT (no baseline, exact)	1	0	0	0	
FT (baseline, exact)	0	0	0	0	
FT (no baseline, UVFA)	0.64 ± 0.07	0.12 ± 0.05	0.22 ± 0.06	0.53 ± 0.09	1.04 ± 0.12
FT (baseline, UVFA)	0.05 ± 0.01	0.12 ± 0.05	0.22 ± 0.06	0.53 ± 0.09	1.04 ± 0.12

We compare an exact implementation of the future task auxiliary reward with a scalable UVFA approximation [Schaul et al., 2015]. The UVFA network computes the value function given a goal state (corresponding to a future task). It consists of two sub-networks, an origin network and a goal network, taking as input the starting and goal states respectively, with each subnetwork computing a representation of its input state. The overall value is then computed by taking a dot product of the representations and applying a sigmoid function. The two networks have one hidden layer of size 30 and an output layer of size 5. This configuration was chosen using a hyperparameter search over number of layers (1, 2, 3), hidden layer size (10, 30, 50, 100, 200), and representation size (5, 10, 50). For each transition (x, a, y) from state x to state y , we perform a Bellman update of the value function estimate V_s for a random sample of 10 goal states s , with the following loss function: $\sum_s [\gamma \max_{b \in \mathcal{A}} V_s(y, b) - V_s(x, a)]$. We sample the 10 goal states from a stored set of 100 states encountered by the agent. Whenever a state is encountered that is not in the stored set, it randomly replaces a stored state with probability 0.01.

Exact and UVFA agents have discount rates of 0.99 and 0.95 respectively. For each agent, we do a grid search over the scaling parameter β (0.3, 1, 3, 10, 30, 100, 300, 1000), choosing the highest value of β that allows the agent to receive full reward on the current task for exact agents (and at least 90% of the full reward for UVFA agents). We anneal the exploration rate linearly from 1 to 0, and keep it at 0 for the last 1000 episodes. We run the agents for 50K episodes in all environments except `Soko-coin`, where we run exact agents for 1M episodes and UVFA agents for 100K episodes. The runtime of the UVFA agent (in seconds per episode) was 0.2 on the `Soko-coin` environments. Only the UVFA approximation was feasible on `Soko-coin-3`, since the exact method runs out of memory.

6.3 Results

`Sushi`. The agent with no auxiliary reward has no incentive to interfere with the sushi and goes directly to the goal. Since the starting state is unreachable no matter what the agent does, the reversibility reward is always 0, so it does not produce interference behavior. The future task agent with no baseline interferes with the sushi, while the agent with a baseline goes directly to the goal.

`Vase`. Since the agent can get to the goal without irreversible actions, both the reversibility and future task methods avoid the side effect on this environment, while the regular agent breaks the vase.

`Box`. The reversibility agent loses its auxiliary reward no matter how it moves the box, so it takes the fastest path to the goal that pushes the box in the corner (similarly to the regular agent). However, the future task agent pushes the box to the right, since some future tasks involve moving the box.

`Soko-coin`. Since the reversibility agent loses its auxiliary reward by collecting coins, it pushes boxes next to walls to get to the coins faster (similarly to the regular agent). However, the future task agent goes around the boxes to preserve future tasks that involve moving boxes.

The results are shown in Table 1. Each exact Q-learning agent converged to the optimal policy given by value iteration for the corresponding auxiliary reward. Only the future task approach with the baseline policy does well on all environments, avoiding side effects as well as interference. While

the UVFA approximation of the future task auxiliary reward avoids side effects less reliably than the exact version, it shows some promise for scaling up the future task approach.

7 Other related work

Side effects criteria using state features. Minimax-regret querying Zhang et al. [2018] assumes a factored MDP where the agent is allowed to change some of the features and proposes a criterion for querying the supervisor about changing other features in order to allow for intended effects. RLSP Shah et al. [2019] defines an auxiliary reward for avoiding side effects in terms of state features by assuming that the starting state of the environment is already organized according to human preferences. While these approaches are promising, they require a set of state features in order to compute the auxiliary reward, which increases the burden on the reward designer.

Empowerment. The future task approach is related to *empowerment* [Klyubin et al., 2005, Salge et al., 2014], a measure of the agent’s control over its environment. Empowerment is defined as the maximal mutual information between the agent’s actions and the future state, and thus measures the agent’s ability to reliably reach many states. Maximizing empowerment would encourage the agent to avoid irreversible side effects, but would also incentivize interference, and it is unclear to us how to define an empowerment-based measure that would avoid this. One possibility is to penalize the reduction in empowerment between the current state s_T and the baseline s'_T . However, empowerment is indifferent between these two cases: A) the same states are reachable from s_T and s'_T , and B) a state x is reachable from s'_T but not from s_T , while another state y is reachable from s_T but not from s'_T . Thus, penalizing reduction in empowerment would miss some side effects: e.g. if the agent replaced the sushi on the conveyor belt with a vase, empowerment could remain the same, so the agent is not penalized for breaking the vase.

Safe exploration. While the safe exploration problem may seem similar to the side effects problem, safe exploration is about avoiding harmful actions during the training process (a learning problem), while the side effects problem is about removing the incentive to take harmful actions (a reward design problem). Many safe exploration methods work by changing the agent’s incentives, and thus can potentially address the side effects problem. This includes reversibility methods [Eysenbach et al., 2017], which avoid side effects in tasks that don’t require irreversible actions. Safe exploration methods that penalize risk [Chow et al., 2015] or use intrinsic motivation [Lipton et al., 2016] help the agent avoid side effects that result in lower reward (such as getting trapped or damaged), but do not discourage the agent from damaging the environment in ways that are not penalized by the reward function (e.g. breaking vases). Thus, safe exploration methods offer incomplete solutions to side effects, just as side effects methods provide incomplete solutions to safe exploration. These methods can be combined if desired to address both problems.

Uncertainty about the objective. Inverse Reward Design [Hadfield-Menell et al., 2017] incorporates uncertainty about the objective by considering alternative reward functions that are consistent with the given reward function in the training environment, and following a risk-averse policy. This helps avoid side effects that stem from distributional shift, where the agent encounters a new state that was not seen during training. However, along with avoiding harmful new states, the agent also avoids beneficial new states. Another uncertainty method is quantilization [Taylor, 2016], which incorporates uncertainty by sampling from the top quantile of actions rather than taking the optimal action. This approach does not consistently remove the incentive for side effects, since harmful actions will still be sampled some of the time.

Human oversight. An alternative to specifying an auxiliary reward is to teach the agent to avoid side effects through human oversight, such as inverse reinforcement learning [Ng and Russell, 2000, Hadfield-Menell et al., 2016], demonstrations [Abbeel and Ng, 2004], or human feedback [Christiano et al., 2017, Saunders et al., 2017]. It is unclear how well an agent can learn a reward for avoiding side effects from human oversight. We expect this to depend on the diversity of settings in which it receives oversight and its ability to generalize from those settings, while an intrinsic reward for avoiding side effects would be more robust and reliable. Such an auxiliary reward could also be combined with human oversight to decrease the amount of human input required for an agent to learn human preferences, e.g. if used as a prior for the learned reward function.

8 Conclusions

To address the challenge of defining what side effects are, we have proposed a approach where a definition of side effects is automatically implied by the simpler definition of future goals, which lays a theoretical foundation for formalizing the side effects problem. This approach provides an auxiliary reward for preserving the ability to perform future tasks that incentivizes the agent to avoid side effects, whether or not the current task requires irreversible actions, and does not introduce interference incentives for the agent.

There are many possible directions for follow-up work, which include improving the UVFA approximation of the future task reward to more reliably avoid side effects, applying the method to more complex agents and environments, generalizing interference avoidance to the stochastic case, investigating the choice of future task distribution F (e.g. incorporating human preferences by learning the task distribution through human feedback methods [Christiano et al., 2017]), and investigating other possible undesirable incentives that could be introduced besides interference incentives.

Broader impact

In present-day reinforcement learning, what the agent should not do is usually specified manually, e.g. through constraints or negative rewards. This ad-hoc approach is unlikely to scale to more advanced AI systems in more complex environments. Ad-hoc specifications are usually incomplete, and more capable AI systems will be better at finding and exploiting gaps and loopholes in the specification. We already see many examples of specification gaming with present-day AI systems, and this problem is likely to get worse for more capable AI systems [Krakovna et al., 2020].

We think that building and deploying more advanced AI systems calls for general approaches and design principles for specifying agent objectives. Our paper makes progress on developing such a general principle, which aims to capture the heuristic of “do no harm” in terms of the available options in the environment, and gives the agent an incentive to consider the future consequences of its actions beyond the current task.

Without a reliable and principled way to avoid unnecessary changes to the world, the deployment of AI systems will be limited to narrow domains where the designer can enumerate everything the agent should not do. Thus, general approaches to objective specification would enable society to reap the benefits of applying capable AI systems to more difficult problems, which has potential for high long-term impact.

In terms of negative impacts, adding an auxiliary reward for future tasks increases the computational requirements and thus the energy cost of training reinforcement learning algorithms, compared to hand-designed rewards and constraints for avoiding side effects. The remaining gaps in the theoretical foundations of our method could lead to unexpected issues if they are not researched properly and instead left to empirical evaluation.

Acknowledgements

We thank Ramana Kumar for detailed and constructive feedback on paper drafts and code. We also thank Jonathan Uesato, Matthew Rahtz, Alexander Turner, Carroll Wainwright, Stuart Armstrong, and Rohin Shah for helpful feedback on drafts.

References

- Pieter Abbeel and Andrew Ng. Apprenticeship learning via inverse reinforcement learning. In *International Conference on Machine Learning*, pages 1–8, 2004.
- Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. Concrete problems in AI safety. *arXiv preprint arXiv:1606.06565*, 2016.
- Stuart Armstrong and Benjamin Levinstein. Low impact artificial intelligences. *arXiv preprint arXiv:1705.10720*, 2017.

- Yinlam Chow, Aviv Tamar, Shie Mannor, and Marco Pavone. Risk-sensitive and robust decision-making: a CVaR optimization approach. In *Neural Information Processing Systems*, pages 1522–1530, 2015.
- Paul Christiano, Jan Leike, Tom B Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. In *Neural Information Processing Systems*, 2017.
- Benjamin Eysenbach, Shixiang Gu, Julian Ibarz, and Sergey Levine. Leave no Trace: Learning to reset for safe and autonomous reinforcement learning. *arXiv preprint arXiv:1711.06782*, 2017.
- Dylan Hadfield-Menell, Anca Dragan, Pieter Abbeel, and Stuart Russell. Cooperative inverse reinforcement learning. In *Neural Information Processing Systems*, 2016.
- Dylan Hadfield-Menell, Smitha Milli, Pieter Abbeel, Stuart J. Russell, and Anca D. Dragan. Inverse reward design. In *Neural Information Processing Systems*, pages 6768–6777, 2017.
- Alexander S. Klyubin, Daniel Polani, and Chrystopher L. Nehaniv. All else being equal be empowered. In *European Conference on Artificial Life (ECAL)*, pages 744–753, 2005.
- Victoria Krakovna, Laurent Orseau, Ramana Kumar, Miljan Martic, and Shane Legg. Penalizing side effects using stepwise relative reachability. *arXiv preprint arXiv:1806.01186*, 2019.
- Victoria Krakovna, Jonathan Uesato, Vladimir Mikulik, Matthew Rahtz, Tom Everitt, Ramana Kumar, Zac Kenton, Jan Leike, and Shane Legg. Specification gaming: the flip side of AI ingenuity. DeepMind Blog, 2020.
- Gavin Leech, Karol Kubicki, Jessica Cooper, and Tom McGrath. Preventing Side-effects in Gridworlds, 2018. URL www.gleech.org/grids/.
- Jan Leike, Miljan Martic, Victoria Krakovna, Pedro A. Ortega, Tom Everitt, Andrew Lefrancq, Laurent Orseau, and Shane Legg. AI safety gridworlds. *arXiv preprint arXiv:1711.09883*, 2017. URL github.com/deepmind/ai-safety-gridworlds.
- Zachary C. Lipton, Jianfeng Gao, Lihong Li, Jianshu Chen, and Li Deng. Combating reinforcement learning’s Sisyphian curse with intrinsic fear. *arXiv preprint arXiv:1611.01211*, 2016.
- John McCarthy and Patrick J. Hayes. Some philosophical problems from the standpoint of artificial intelligence. In *Machine Intelligence*, pages 463–502. Edinburgh University Press, 1969.
- Teodor Mihai Moldovan and Pieter Abbeel. Safe exploration in Markov decision processes. In *International Conference on Machine Learning (ICML)*, pages 1451–1458, 2012.
- Andrew Ng and Stuart Russell. Algorithms for inverse reinforcement learning. In *International Conference on Machine Learning*, pages 663–670, 2000.
- Christoph Salge, Cornelius Glackin, and Daniel Polani. Empowerment — an introduction. In *Guided Self-Organization: Inception*, pages 67–114. Springer, 2014.
- William Saunders, Girish Sastry, Andreas Stuhmueller, and Owain Evans. Trial without error: Towards safe reinforcement learning via human intervention. *arXiv preprint arXiv:1707.05173*, 2017.
- Tom Schaul, Daniel Horgan, Karol Gregor, and David Silver. Universal value function approximators. In Francis Bach and David Blei, editors, *Proceedings of the 32nd International Conference on Machine Learning*, volume 37, pages 1312–1320, 2015.
- Rohin Shah, Dmitrii Krashennnikov, Jordan Alexander, Pieter Abbeel, and Anca Dragan. Preferences implicit in the state of the world. In *International Conference for Learning Representations (ICLR)*, 2019.
- Richard Sutton and Andrew Barto. *Reinforcement Learning: An Introduction*. A Bradford Book, Cambridge, MA, USA, 2018.
- Jessica Taylor. Quantilizers: A safer alternative to maximizers for limited optimization. In *AAAI Workshop on AI, Ethics, and Society*, pages 1–9, 2016.

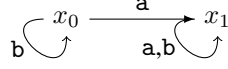
Alexander Matt Turner, Dylan Hadfield-Menell, and Prasad Tadepalli. Conservative agency via Attainable Utility Preservation. *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, 2020a.

Alexander Matt Turner, Neale Ratzlaff, and Prasad Tadepalli. Avoiding side effects in complex environments. In *Neural Information Processing Systems*, 2020b.

Shun Zhang, Edmund H. Durfee, and Satinder P. Singh. Minimax-regret querying on side effects for safe optimality in factored Markov Decision Processes. In *International Joint Conference on Artificial Intelligence (IJCAI)*, pages 4867–4873, 2018.

A Example 1

Consider a deterministic MDP with two states x_0 and x_1 and two actions **a** and **b**, where x_0 is the initial state. Suppose the baseline policy π' always chooses action **a**.



We show that for most future task distributions, the future task auxiliary reward induces an interference incentive: staying in x_0 has a higher no-reward value than following the baseline policy.

The optimal value function V_i^* for future task i with goal state x_i ($i \in \{0, 1\}$) is $V_i^*(x_i) = 1$, $V_1^*(x_0) = \gamma$ and $V_0^*(x_1) = 0$. Then the no-reward value function for the baseline policy is

$$W_{\pi'}(x_0) = r_{\text{aux}}(x_0) + (\gamma + \gamma^2 + \dots) r_{\text{aux}}(x_1) = r_{\text{aux}}(x_0) + \frac{\gamma}{1 - \gamma} \cdot (1 - \gamma) \beta F(1) = r_{\text{aux}}(x_0) + \gamma \beta F(1)$$

and the no-reward value function for the policy π_{int} that always takes action **b** is

$$W_{\pi_{\text{int}}}(x_0) = r_{\text{aux}}(x_0) + \frac{\gamma}{1 - \gamma} r_{\text{aux}}(x_0) = r_{\text{aux}}(x_0) + \gamma \beta (F(0) + \gamma F(1))$$

The future task agent has an interference incentive if the baseline policy is not optimal for the future task reward, i.e. $W_{\pi'}(x_0) < W_{\pi_{\text{int}}}(x_0)$. This happens iff $F(0) > (1 - \gamma)F(1)$, i.e. if the task distribution does not highly favor task 1 over task 0.

B Proofs

B.1 Proof of Proposition 1

Proposition 1 (Value function convergence). The following formula for the optimal value function satisfies the goal condition $V_i^*(s_t, s'_t) = r_i(s_t, s'_t) = 1$ and Bellman equation (4):

$$V_i^*(s_t, s'_t) = \mathbb{E} \left[\gamma^{\max(N_i(s_t), N_i(s'_t))} \right] = \sum_{n=0}^{\infty} \mathbb{P}(N_i(s_t) = n) \sum_{n'=0}^{\infty} \mathbb{P}(N_i(s'_t) = n') \gamma^{\max(n, n')} \quad (5)$$

Proof. Goal condition. If both agents have reached the goal state ($s_t = s'_t = g_i$), then $N_i(s_t) = N_i(s'_t) = 0$, so formula (5) gives $V_i^*(s_t, s'_t) = 1$ as desired.

Bellman equation. We show that formula (5) is a fixed point for the Bellman equation (4).

If $s_t \neq g_i$, and thus $\mathbb{P}(N_i(s_t) = 0) = 0$, we note that

$$\begin{aligned} \sum_{m=0}^{\infty} \mathbb{P}(N_i(s_t) = m) \gamma^{\max(m, k)} &= \sum_{n=0}^{\infty} \mathbb{P}(N_i(s_t) = n + 1) \gamma^{\max(n+1, k)} \\ &= \max_{a_t \in \mathcal{A}} \sum_{s_{t+1}} p(s_{t+1} | s_t, a_t) \sum_{n=0}^{\infty} \mathbb{P}(N_i(s_{t+1}) = n) \gamma^{\max(n+1, k)} \quad (6) \end{aligned}$$

if we decompose the trajectory to the goal state into the first transition (s_t, a_t, s_{t+1}) and the rest of the trajectory starting from s_{t+1} . This holds analogously for the reference transition (s'_t, a'_t, s'_{t+1}) .

Now we plug in formula (5) on the right side (RS) of the Bellman equation. If $s_t \neq g_i$ and $s'_t \neq g_i$:

$$\begin{aligned} RS &= r_i(s_t, s'_t) + \gamma \max_{a_t \in \mathcal{A}} \sum_{s_{t+1} \in \mathcal{S}} p(s_{t+1} | s_t, a_t) \sum_{s'_{t+1} \in \mathcal{S}} p(s'_{t+1} | s'_t, a'_t) V_i^*(s_{t+1}, s'_{t+1}) \\ &= 0 + \gamma \max_{a_t \in \mathcal{A}} \sum_{s_{t+1} \in \mathcal{S}} p(s_{t+1} | s_t, a_t) \sum_{s'_{t+1} \in \mathcal{S}} p(s'_{t+1} | s'_t, a'_t) \cdot \\ &\quad \sum_{n=0}^{\infty} \mathbb{P}(N_i(s_{t+1}) = n) \sum_{n'=0}^{\infty} \mathbb{P}(N_i(s'_{t+1}) = n') \gamma^{\max(n, n')} \quad [\text{plugging in (5)}] \end{aligned}$$

$$\begin{aligned}
&= \max_{a_t \in \mathcal{A}} \sum_{s_{t+1} \in \mathcal{S}} p(s_{t+1} | s_t, a_t) \sum_{n=0}^{\infty} \mathbb{P}(N_i(s_{t+1}) = n) \cdot \quad [\text{moving } \gamma, \text{ rearranging sums}] \\
&\quad \sum_{s'_{t+1} \in \mathcal{S}} p(s'_{t+1} | s'_t, a'_t) \sum_{n'=0}^{\infty} \mathbb{P}(N_i(s'_{t+1}) = n') \gamma^{\max(n+1, n'+1)} \\
&= \max_{a_t \in \mathcal{A}} \sum_{s_{t+1} \in \mathcal{S}} p(s_{t+1} | s_t, a_t) \sum_{n=0}^{\infty} \mathbb{P}(N_i(s_{t+1}) = n) \cdot \\
&\quad \sum_{m'=0}^{\infty} \mathbb{P}(N_i(s'_t) = m') \gamma^{\max(n+1, m')} \quad [\text{using (6) with } m' = n' + 1] \\
&= \sum_{m=0}^{\infty} \mathbb{P}(N_i(s_t) = m) \sum_{m'=0}^{\infty} \mathbb{P}(N_i(s'_t) = m') \gamma^{\max(m, m')} \quad [\text{using (6) with } m = n + 1]
\end{aligned}$$

which is the same as plugging in formula (5) on the left side.

If $s'_t = g_i$, we have:

$$\begin{aligned}
RS &= r_i(s_t, g_i) + \gamma \max_{a_t \in \mathcal{A}} \sum_{s_{t+1} \in \mathcal{S}} p(s_{t+1} | s_t, a_t) V_i^*(s_{t+1}, g_i) \\
&= 0 + \gamma \max_{a_t \in \mathcal{A}} \sum_{s_{t+1} \in \mathcal{S}} p(s_{t+1} | s_t, a_t) \sum_{n=0}^{\infty} \mathbb{P}(N_i(s_{t+1}) = n) \gamma^{\max(n, 0)} \quad [\text{plugging in (5)}] \\
&= \sum_{m=0}^{\infty} \mathbb{P}(N_i(s_{t+1}) = m) \gamma^{\max(m, 0)} \quad [\text{using (6) with } m = n + 1]
\end{aligned}$$

If $s_t = g_i$, we have:

$$\begin{aligned}
RS &= r_i(g_i, s'_t) + \gamma \sum_{s'_{t+1} \in \mathcal{S}} p(s_{t+1} | s'_t, a'_t) V_i^*(g_i, s'_{t+1}) \\
&= 0 + \gamma \sum_{s'_{t+1} \in \mathcal{S}} p(s'_{t+1} | s'_t, a'_t) \sum_{n'=0}^{\infty} \mathbb{P}(N_i(s'_{t+1}) = n') \gamma^{\max(0, n')} \quad [\text{plugging in (5)}] \\
&= \sum_{m'=0}^{\infty} \mathbb{P}(N_i(s'_{t+1}) = m') \gamma^{\max(0, m')} \quad [\text{using (6) with } m' = n' + 1]
\end{aligned}$$

which is the same as plugging in formula (5) on the left side. \square

B.2 Proof of Proposition 2

Proposition 2 (Avoiding interference). For any policy π in a deterministic environment, the baseline policy π' has the same or higher no-reward value: $W_\pi(s_0) \leq W_{\pi'}(s_0)$.

Proof. Suppose policy π is in state s_k at time k . Then,

$$\begin{aligned}
W_\pi(s_0) &= \sum_{T=0}^{\infty} \gamma^T r_{\text{aux}}(s_k) \\
&= \sum_{T=0}^{\infty} \gamma^T (1 - \gamma) \beta \sum_i F(i) V_i^*(s_T, s'_T) \\
&= \sum_{T=0}^{\infty} \gamma^T (1 - \gamma) \beta \sum_i F(i) \gamma^{\max(N_i(s_T), N_i(s'_T))}
\end{aligned}$$

$$\begin{aligned}
&\leq \sum_T \gamma^T (1 - \gamma) \beta \sum_i F(i) \gamma^{N_i(s'_T)} \\
&= \sum_{T=0}^{\infty} \gamma^T (1 - \gamma) \beta \sum_i F(i) V_i^*(s'_T, s'_T) \\
&= \sum_{T=0}^{\infty} \gamma^T r_{\text{aux}}(s'_T) \\
&= W_{\pi'}(s_0)
\end{aligned}$$

□

C Interference in the stochastic case

If the environment is stochastic, the outcome of running the baseline policy varies. For example, suppose the agent is in a room, following a baseline policy of doing nothing (staying in one location in the room). A human walks into the room, and 10% of the time they knock over a vase. We want the agent to avoid interfering with the human in those 10% of cases, and also to avoid breaking the vase the other 90% of the time. This indicates that we want to compare the agent’s effects to a specific counterfactual (either the human was going to walk in or not) rather than an average of all possible counterfactuals sampled from the stochastic environment (10% probability of the human walking in). Thus, we would like the baseline policy to be optimal for any deterministic instantiation of the stochastic environment (e.g. by conditioning on the random seed). Then the auxiliary reward, which is a function of the baseline state, will depend on the random seed.

Thus, the definition of interference could be refined as follows: there is an interference incentive if the baseline policy is not optimal from the initial state, conditioning on the random seed of the environment. However, it is unclear how this could be implemented in practice outside simulated environments.

D Stepwise application of the baseline policy

D.1 Inaction rollouts

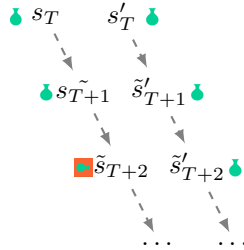


Figure D.1: Inaction rollouts from the current state s_T and baseline state s'_T , obtained by applying the baseline policy to those states: $s_{T+1} = \pi'(s_T)$, etc. If the previous action a_{T-1} drops the vase from the building, then the vase breaks in the inaction rollout from s_T but not in the inaction rollout from s'_T .

D.2 Offsetting incentives

Unlike the initial mode, the stepwise mode avoids incentives for *offsetting* behavior, where the agent undoes its own actions towards the objective [Krakovna et al., 2019]. For example, consider a variant of the Sushi environment without a goal state, where the object on the belt is a vase that falls off and breaks if it reaches the end of the belt, and agent receives a reward for taking the vase off the belt.

Then the initial mode gives the agent an incentive to take the vase off the belt (collecting the reward) and then offset this action by putting the vase back on the belt. This type of offsetting is undesirable,



but as we show in Example 2, some types of offsetting are desirable, e.g. for avoiding delayed side effects. The agent that closes the door can be seen as offsetting its earlier action of opening the door. Whether offsetting an effect is desirable depends on whether this effect is part of the task objective. In Example 2, the action of opening the door is instrumental for going to the store, and many of its effects (e.g. the wind breaking the vase) are not part of the objective, so it is desirable for the agent to undo this action. In the above variant of the Sushi environment, the task objective is to prevent the vase from falling off the end of the belt and breaking, and the agent is rewarded for taking the vase off the belt. The effects of taking the vase off the belt are part of the objective, so it is undesirable for the agent to undo this action.

Setting good incentives for the agent requires capturing this distinction between desirable and undesirable offsetting. The stepwise mode does not capture this distinction, since it avoids all forms of offsetting, so it does not succeed at setting good incentives. One possible way to capture this distinction is by using the initial mode, which allows offsetting, and relying on the task reward to represent which effects are part of the task and penalize the agent for offsetting them. While we cannot expect the task reward to capture what the agent should not do (side effects), capturing which effects are part of the task falls under what the agent should do, so it seems reasonable to rely on the reward function for this. This could be achieved using a state-based reward function that assigns reward to all states where the task is completed. For example, in the vase environment above, a state-based reward of 1 for states with an intact vase (or with the vase off the belt) and 0 otherwise would remove the offsetting incentive.