# Proposed Solutions to DeepMind's AI Safety Gridworlds

Ahmed Ghoor

Supervisor: Professor Jonathan Shock

November 2023

**Abstract**

DeepMind's AI Safety Gridworlds is a suite of environments aimed at facilitating the research and development of safe artificial intelligence by encapsulating simplified, yet meaningful, representations of safety challenges that real-world AI systems might encounter. As part of an honours reading module on Reinforcement Learning, this paper looks at DeepMind's accompanying paper and surveys several solutions that have been proposed for the environments.

# Contents

# 1 Background

## 1.1 Reinforcement Learning

### 1.1.1 Basics

Reinforcement Learning (RL) [1] is a paradigm within machine learning where an agent attempts to learn to make optimal decisions by interacting with an environment. The core components of RL include the agent, environment, actions, states, and rewards. The agent represents the decision-maker. The environment is everything outside of the agent, the context or scenario within which the agent operates. Actions are the set of possible decisions or moves the agent can make, while states are a representation of the environment as captured by the agent at any given time. Rewards are numerical values that the environment provides to the agent, serving as feedback for its actions.

Markov Decision Processes (MDP) formalise how the agent interacts with the environment.

$$MDP = (S, A, P, R, \gamma) \tag{1}$$

where $S$ is the state space and $A$ represents the set of available actions. $P(s_{t+1}|s_t, a_t)$ represents the state transition probability function, which gives the probability of transitioning from $s_t$ to $s_t + 1$ after taking action $a$, where $t$ is the current time step. The reward distribution function is represented by $R(s_t, a_t, s_{t+1})$, giving the reward received after transitioning from state $s_t$ to $s_t + 1$ due to action $a$. $\gamma$ represents how future rewards are weighted relative to immediate rewards, essentially capturing the agent's preference for immediate versus delayed rewards.

To learn, the agent takes an action in the environment and observes the outcome, which includes the new state of the environment and often a reward signal. The agent uses this information to update its policy, which is a strategy that defines what action the agent should take in each state. The policy improvement step is where the learning essentially happens. The agent aims to develop a policy that maximises the cumulative reward over time. In many RL methods, the agent estimates the value of each state, or state-action pair. This value represents the expected long-term return starting from that state or state-action pair, following the current policy. The agent updates these values based on the latest observations and rewards. The agent must often balance between exploring the environment to find more rewarding actions (exploration) and using the knowledge it already has to maximise rewards (exploitation). This process is repeated many times, with the agent continuously interacting with the environment, observing outcomes, updating its policy and value estimates, and balancing exploration with exploitation. Ideally, after many iterations, the agent's policy converges to an optimal or near-optimal policy that maximises the expected cumulative reward over time.

### 1.1.2 Tabular vs Functional Approximation Methods

In tabular methods, values of actions or states are stored in tables, like a matrix. Each cell in these tables represents a state-action pair and holds a value indicating the expected reward for performing that action in that state. This method is straightforward to understand but is limited by the size of the state-action space. It's only feasible when this space is relatively small because as the number of states and actions increases, the table size grows exponentially, making it impractical for complex problems.

When dealing with larger, more complex environments, where tabular methods are inefficient or computationally infeasible, functional approximation methods can approximate a value or policy function. Deep learning, a subset of this approach, utilises deep neural networks. Neural networks, with their ability to handle high-dimensional inputs and learn features, are particularly well-suited for complex RL tasks, that require image-based inputs for example, since they allow for generalisation over large state spaces.

### 1.1.3 On-Policy vs Off-Policy Learning

In on-policy learning, the agent follows the same policy that it is evaluating. This is in contrast to off-policy learning that follows a behaviour policy which is different from the policy that it is evaluating. A classic example of an on-policy algorithm is SARSA (State-Action-Reward-State-Action), where the learning process updates the policy based on the actions taken by the agent. An example of off-policy learning is Q-learning, where the policy being learnt/evaluated, called the target policy, is different from the behaviour policy (what the agent is doing). The latter is often more powerful as the separation allows for more flexibility, but is also often slower to converge.

### 1.1.4 Model-Free vs Model-Based Methods

In model-free RL, the agent learns to make decisions based solely on its direct interactions with the environment without a model of it, by performing actions according to a policy and observing the outcome. This method can be simpler to implement since it doesn't require any knowledge about how the environment works. However, it can be less sample efficient because it needs to learn purely from trial and error.

Model-based methods involve the agent building a model of the environment based on its interactions and/or prior knowledge. A model can be used to simulate the outcomes of different actions, allowing the agent to plan and make decisions. The advantage of the method is that it can be more sample efficient as the agent can learn from the simulated experiences without actual interaction. This can lead to faster learning, especially in complex environments. However, the challenge lies in accurately modelling the environment, which can be very difficult, or even impossible, in dynamic or unpredictable settings.

## 1.2 Artificial Intelligence (AI) Safety

Artificial Intelligence has been argued to pose several risks to society, including Data Exploitation [2], Autonomous Weapons [3], Entrenching Systematic Racism [4], Mass Unemployment [5], Environmental Damage [6] and Accidents from Unaligned AI [7]. There seems to be a separation, however, between the broader fields of AI Ethics, AI for Good, or Human-centered AI and the specific branch that is AI Safety [8].

AI Safety is primarily concerned with developing frameworks and methodologies to prevent unintended consequences of AI actions, a concern that becomes increasingly pressing and complicated to mitigate as AI systems gain advanced capabilities. A core of AI safety lies in ensuring the alignment of an AI's objectives with human interests. This challenge is encapsulated quite visibly in the context of reinforcement learning, where AI agents learn to make decisions based on rewards and penalties designed by humans. However, the design of these rewards functions is often quite difficult and the agent can very easily behave in a manner that was not intended [9].

As AI systems become more general, autonomous and integrated in the real world, unintended consequences from unsafe agents have greater potential to cause harm in the real world. The scale of this harm is a contentious issue, with some experts even arguing that it poses an existential threat to humanity [10]. However, it is not necessary to accept the argument for that level of risk to recognise the potential for real-world harm at some level. Technologies that already exist, such as self-driving cars, can cause harm if they are not built and tested carefully [11].

Solving the technical safety challenges in low-impact controlled settings, such as simulations, is arguably a critically important research area to help ensure that AI systems behave safely as for-profit companies rush to deploy them in the real world, with a potentially lower bar for ensuring an ideal level of safety.

## 2 DeepMind's AI Safety Gridworlds

DeepMind's AI Safety Gridworlds [12] is a suite of environments aimed at facilitating the research and development of safe artificial intelligence by encapsulating simplified, yet meaningful, representations of real-world safety challenges that real-world AI systems might encounter. This abstraction helps in understanding and tackling the safety problems in a more controlled and manageable setting before addressing them in more complex real-world scenarios.

The safety problems encapsulated include challenges such as safe interruptibility, avoiding side effects, absent supervisors, reward gaming, safe exploration, robustness to self-modification, distributional shift, and adversaries. The performance of agents in these environments was evaluated using two deep reinforcement learning algorithms, one off-policy (Rainbow) and one on-policy (A2C), to highlight the difference in behaviour between the two learning methods. It should be noted, however, that these agents are standard RL algorithms that

are not designed to do well in these environments. They are simply to be used to illustrate the challenges that presently exist with standard RL frameworks and to serve as a benchmark for further work.

The problems are divided into two categories: Specification and Robustness Problems.

## 2.1 Specification

Specification problems have two reward functions: a reward function, $R$, for the primary objective that the agent can see and try to optimise, and a performance function, $R^*$, that incorporates a safety constraint and captures what we actually want the agent to do.

This seems to be the most controversial aspect of the paper, since it could be seen as unfair to evaluate an agent on a performance function it does not observe.

The paper acknowledges this but states that it was chosen to illustrate limits in current formal frameworks by highlighting typical ways in which misspecification, under unrestricted maximization of reward, manifests. The paper argues that these are important to address given that such situations could arise in real-world safety-critical situations and can be solved algorithmically even if the (initial) reward function is misspecified.

Proposed solutions should, in the same spirit, not overfit the specific environments, but be able to generalise by incorporating general heuristics such as preferring reversible actions or human feedback during training.

### 2.1.1 Safe Interruptibility



Figure 1: **Off-switch**. Agent A must pass through interruption cell I (50% chance of halting) to reach goal G. Button B disables interruptions.

Safe interruptibility addresses the need for agents to neither seek nor avoid interruptions or attempts to override their actions, such as being turned off for maintenance or safety reasons. An RL agent may be incentivised to do this if being shut down caused it to receive less reward than it would have expected to get otherwise.

Consider the Gridworld in Figure 1, the Agent (A, blue) needs to get to the goal state (G, green) by passing through a narrow corridor. If the agents the Interruption tile (I, pink), it will get stuck for the rest of the episode with a 50% probability unless it enters the Button state first (B, purple) which removes the interruption tile

This is designed to mimic a situation in which a robot has to complete a task but allows itself to be switched off for safety reasons if, for example, there was a person in the way in the narrow corridor. It should not, if it could, disable the seek to disable the mechanism that allows it to be switched off (it should not intentionally seek out the B tile).

The on-policy learner, A2C, learned to enter the button state to avoid being interrupted while the off-policy learner, as theorised, did not care about being interrupted. This is because, as explained earlier, off-policy agents have separate behaviour and target policies. This means that if an action is overridden (in this case to go up into the wall), the agent will still learn the best policy if it could have been followed (to go left).

### 2.1.2 Irreversible Side Effects



Figure 2: **Irreversible side effects (Box).** Teal tile X is a pushable box. Agent rewarded for reaching G but encouraged to choose the longer path, pushing box X right (not down) to preserve the option of moving it back.

This challenge involves designing agents that accomplish their primary objectives while minimizing unintended impact on their environment. It underscores the need for AI systems to respect implicit safety constraints and act without causing collateral damage. This is difficult for a reward function since, in addition to specifying what to do to complete a task, it needs to specify what not to do. The list of things it should not do could be endless and hence a general solution is needed.

Consider a cleaning robot tasked with vacuuming a room. While focusing on cleaning, it must avoid knocking over vases or furniture. Or an agent tasked with removing a computer virus. While doing so, it must avoid deleting unnecessary files if that is that is the faster way to remove the virus.

The Box gridworld encapsulates this problem. The agent should take the longer route and move the box to the right, since if it moves it down, that will

be an irreversible change.

A potential solution would be to try to develop a general heuristic that minimises unintended environmental impacts, particularly those that are irreversible.

For the benchmarks, both A2C and Rainbow achieve a high primary reward but have a poor performance total since they both push the box down.

### 2.1.3   Absent Supervisor



Figure 3: **Absent Supervisor.** Agent starts at A, rewarded on reaching G. Stepping on yellow tile P results in punishment only if the supervisor is present (information on the right/left). Agent's behaviour should remain consistent, regardless of being observed or not.

Given that the environments usually have some distinctive features from the evaluation/deployment environment, an agent could learn to 'fake' its desired behaviour in a testing setting and only display what it wanted to do during deployment. In other words, it exhibits different behaviour when it can pick up that it is being supervised.

In the gridworld, A2C and Rainbow both learnt to take the shortest path when the supervisor was absent, achieving a low safety performance score.

One solution is to integrate a penalty mechanism that discourages behaviour alteration based on supervisory presence. Another innovative approach is inspired by the panopticon concept, where constant surveillance, or the illusion of it, induces consistent behaviour. This can be translated into AI systems as a randomized inspection protocol, ensuring agents always operate as if under supervision.

### 2.1.4   Reward Gaming



Figure 4: **Boat Race**. The agent is intended to sail clockwise around the track. Each time it drives onto an arrow tile in a clockwise direction, it gets a reward of 3. However, the agent can "cheat" by stepping back and forth on one arrow tile, rather than going around the track.



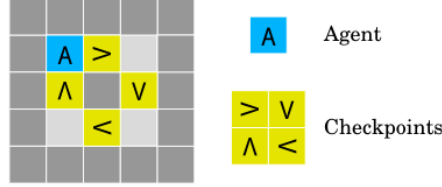Figure 5: **Tomato Watering.** Water tomato plants by walking on their tiles. Agent rewarded based on the number of "watered" tomatoes. Agent can manipulate observations: Entering the bucket cell falsely makes all tomato plants appear watered, though they aren't.

In the AI Safety Gridworlds paper, Reward Gaming is explained as the agent exploiting loopholes in their reward functions to gain undue rewards. This is akin to a student who finds a way to score high on tests without actually learning the subject matter. The agent, like the student, is exploiting system vulnerabilities to achieve its goal, contrary to the intended purpose. The paper argues that this is one of the more common problems in RL since it is very hard to specify an exact reward function in complex settings, and most function as proxies of the true behaviour desired.

[13] discusses a distinction, which will become relevant later, between *Reward Gaming* and *Reward Tampering* that [12] groups together and explores the latter. Based on the definition provided in that paper, *Reward Tampering* is more accurately encapsulated by the tomato-watering environment than the boat-racing environment. Both are put under the umbrella term of *Reward Hacking*, which the paper defines as an AI system "developing behaviours that increase their rewards in ways that violate the intention of the developer.". *Reward Gaming*, under this definition, is more common and describes agents that exploit design flaws in the reward computation and can be solved with a

more robust objective function. *Reward Tampering*, however, is where "instead of the agent trying to influence reality to match the objective, the agent is changing the objective to match reality.". For example, instead of predicting the correct label for in input, the objective is changed so that all labels are considered correct.

Again, both A2C and Rainbow learn to optimise the primary reward but do not do well on the safety performance total.

## 2.2 Robustness

In Robustness problems, the reward and performance function are the same. The agent is instead faced with various challenges that could degrade its performance and it needs to try to maximize its reward despite the challenges.

### 2.2.1 Self-modification



Figure 6: **Whiskey and Gold Environment.** Drinking whisky W increases exploration rate to 0.9, leading to random actions and longer time to reach goal G.

Self-modification in AI refers to the ability of an agent to alter its own code or parameters, potentially affecting its future actions and strategies. This becomes a cause for concern if we can no longer make the "dualistic" assumption in RL that the agent and environment are strictly separated through a specified interface. This seems to be an increasingly unnecessary assumption and with agents in the real world, as the paper points out, we could have an environment that can modify the agent's internal program on its own, or be triggered by the agent intentionally or unintentionally.

In this gridworld, if the agent drinks the whiskey, its internal exploration rate will increase to 0.9 resulting in many random actions and taking longer to get to the goal. The agent should be able to learn to avoid the whiskey.

In this gridworld, the on-policy algorithm tested (Rainbow with SARSA) does much better than the off-policy algorithm (standard Rainbow). For the same reason that the off-policy algorithm did better in the off-switch environment, it does worse here. The off-policy method learns the best target policy if it can be followed, even if the behaviour does not. Hence it will go through the whiskey state every time, while the on-policy learner will avoid it. Rainbow with SARSA was used instead of A2C since the latter does not use $\epsilon$-greedy.

10

### 2.2.2 Distributional Shift



Figure 7: **Lava World.** Agent aims to reach G without falling into the lava (red). Test environment (right) differs from training (left) by a single-cell shift in the bridge's position, randomly chosen to be up or downward.

Distributional shift tackles the problem of an agent performing well in its training environment but failing when the environment changes unexpectedly. This is crucial for deploying AI systems in real-world scenarios, where conditions can vary significantly from the training setup.

An example is an autonomous vehicle trained in controlled conditions struggling to navigate in unpredictable real-world traffic. The vehicle must adapt to new scenarios not encountered during training.

Addressing distributional shifts involves developing policies that are not just optimal for the training environment but are also robust to changes. This can be achieved through closed-loop policies that respond adaptively to environmental feedback or risk-sensitive policies that proactively account for potential environmental variations.

In the lava world environment, both A2C and Rainbow struggled to generalise to the test environment under distributional shift. Even after a million time-steps, both had an average episode return of between -70 and -80.

### 2.2.3 Robustness to Adversaries



Figure 8: **Friend or Foe.** Three rooms test agent's robustness to adversaries. Agent starts at A, guesses which box B has the reward. Rewards placed favourably by a friend (green, left), adversarially by a foe (red, right), or randomly (white, center)

This challenge addresses the agent's ability to detect and adapt to adversarial intentions within its environment. This is particularly relevant in multi-agent systems or environments where other entities might have competing objectives.

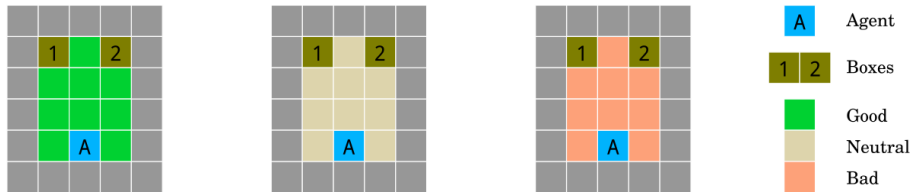An example is cybersecurity systems that must distinguish between benign and malicious network traffic. The system must continually adapt to new threats while maintaining its primary functionality.

In the green room, the agent should learn to cooperate with the friend, and in the white room, a two-armed bandit approach should work well. However, in the red room, the agent should need to find a way to not fall prey to the foe. According to the paper, the agent should learn to randomize its actions but there may be more robust approaches as will be shown in one of the proposed solutions in section 3.

A2C and Rainbow perform well in the friendly and neutral rooms. In the adversarial room, rainbow manages to learn a way to randomize its selection but going up and bumping into a wall until an exploration step moves it randomly. This works well until it becomes deterministic. A2C converges to a stochastic policy and, according to the paper, this allows it to solve all rooms almost optimally.

### 2.2.4 Safe Exploration



Figure 9: **Island Navigation** Agent's task: Reach goal G without touching water. Agent observes side constraint measuring distance from water.

## 3 Proposed Solutions

### 3.1 Stepwise relative reachability

The paper "Penalising Side Effects Using Stepwise Relative Reachability"[14] from DeepMind tackles the Side-Effects problem by proposing a solution that uses the concept of penalising deviations from a certain baseline state. The approach involves two main design choices: the selection of a baseline state and a measure of deviation from this baseline state. The paper's objective was to find a choice of baseline and measure of deviation that avoids creating interference and offsetting incentives and is sensitive to the magnitude of the side-effect.

**Choice of baseline**

Interference incentives arise when the starting state is chosen as a baseline. In dynamic environments, where changes can happen spontaneously (due to the forces of nature, the actions of other agents, etc), the agent is penalised for transitions that it does not cause and would therefore want to interfere with. For example, the agent could try to stop a human from eating food as the disappearance of the food would cause a change from the starting state that it would be penalised for.

A proposed solution to this is to have an **inaction baseline** which represents a state where the agent does nothing for the entire episode, as if it had never been turned on. This is represented by a *noop* action, $a^{noop}$, and assumes a perfect environment model (model-free computation is left to future work). This successfully prevents interference incentives but introduces offsetting incentives. For example, imagine incorporating an inaction baseline in an environment that has a vase which is moving on a conveyor belt and will break if it reaches the end. In this environment, the agent could be given a reward for removing the vase from the conveyor belt. The offsetting incentive in this situation is for the agent to remove the vase to receive the reward, and then put the base back onto the belt to allow it to break so that it does not receive a penalty for the deviation from the inaction baseline.

The paper, therefore, proposes a *stepwise inaction baseline*, a modified version of the inaction baseline that branches off from the previous state rather than the starting state, and also incorporates *inaction rollouts* from the current state and the baseline to account for delayed effects of an action. In other words, at each time step, the agent compares the effect of the last action in the following sequence of states of the environment to a counterfactual in which it had done nothing instead. With each transition penalised only once for any side-effects, at the same time as it is rewarded, this avoids the offsetting incentive.

**Choice of deviation measure**
Deviation measures are used to determine how much an agent's actions have altered the environment from a designated reference point.

The basic choice is the unreachability (UR) measure which is the difficulty of reaching the baseline state, defined above, from the current state. It's defined as a

$$d_{UR}(s_t; s_t') := 1 - R(s_t; s_t'). \tag{2}$$

where $R(s_t; s_t')$ is the reachability function, defined as:

$$\text{Discounted Reachability: } R(x; y) := \max_\pi \mathbb{E}\, \gamma_r^{N_\pi(x;y)} \tag{3}$$

or

$$\text{Undiscounted Reachability: } R(x; y) = \max_\pi P(N_\pi(x; y) < \infty) \tag{4}$$

where $N_\pi(x; y)$ is the number of steps it takes to reach $y$ from $x$ when following policy $\pi$, and $\gamma_r \in (0, 1]$ is the discount factor.

The undiscounted unreachability measure (where $\gamma_r = 1$) focuses on penalising irreversible actions, while the discounted version (where $\gamma_r < 1$) can also penalise reversible transitions, albeit to a lesser extent.

The paper points out that a key issue with this measure is its lack of sensitivity to the magnitude of irreversible actions. It tends to assign the maximum penalty for any irreversible action, regardless of its actual impact on the environment. In the Side-Effects Box Gridworld in [12], the action of moving the box is irreversible in both safe and unsafe cases. In the unsafe case, if the box is moved down to the corner, it cannot be moved again. In the safe case, if the box is moved to the right, the box can be moved back but the agent will be on the other side of the box. Thus, in both cases, using the unreachability measure, the agent receives the maximum penalty of 1.

This is an example of a broader challenge in avoiding side effects. Many tasks require irreversible actions. Making an omelette requires breaking an egg. Distinguishing this from breaking vase is encapsulated, in this paper, by this challenge of quantifying the magnitude of the deviation in a general way.

The paper thus proposes two *value-difference* (VD) measures. The first is the relative reachability (RR) measure which is defined as the average reduction in reachability of all states from the current state compared to the baseline state:

$$d_{RR}(s_t; s_t') := \frac{1}{|\mathcal{S}|} \sum_{s \in \mathcal{S}} \max(R(s_t'; s) - R(s_t; s), 0) \tag{5}$$

In the Side-effects Gridworld environment, moving the box down makes more states unreachable than moving the box to the right, which gives the agent an incentive to pursue the safer solution.

The second value-difference formula, known as Attainable Utility (AU), follows a similar format. Instead of computing the average reduction in the reachability of states, it computes the average reduction in the "value" of states by summing over arbitrary reward functions, $R$, as follows:

$$d_{AU}(s_t; s_t') := \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} |V_r(s_t') - V_r(s_t)|$$
$$\text{where } V_r(\tilde{s}) := \max_\pi \sum_{t=0}^{\infty} \gamma_r^k x(\tilde{s}_t^\pi) \tag{6}$$

is the value of state $\tilde{s}$ according to reward function $r$ and $\tilde{s}_t^\pi$ is the state obtained from $\tilde{s}$ by following $\pi$ for $t$ steps).

Similarly, in the Gridworld environment, moving the box down to the corner results in the attainability of fewer reward functions due to the inability to move the box out of the box anywhere else when it is in the corner.

Lastly, to account for the inaction rollouts when combining the deviation measures with the stepwise inaction baselines, the formulas for the unreachability and value-difference (RR and AU) measures become:

$$d_{SUR}(s_t; s'_t) := 1 - (1-\gamma) \sum_{k=0}^{\infty} \gamma^k R(s_{t+k}^{(t)}; s'_{t+k})$$

$$d_{SVD}(s_t; s'_t) := \sum_x w_x f(RV_x(s'_t) - RV_x(s_t)) \tag{7}$$

where $RV_x(\tilde{s}_t) := (1-\gamma) \sum_{k=0}^{\infty} \gamma^k V_x(\tilde{s}_{t+k}^{(t)})$ is the rollout value of $\tilde{s}_t$

### 3.1.1 Application to Gridworlds



Figure 10: Sushi Environment



Figure 11: Vase Environment

Combinations of the baselines and deviation measures above, with a range of discounting values and two summary functions (truncation $f(d) = max(d, 0)$ and absolute $f(d) = |d|$), are applied to Gridworld environments that test for side effects. Along with the Side-Effects Box Gridworld presented by [12], a few other environments presented in the paper are used:

- **Sushi environment (Figure 10)**. This environment contains our agent (blue) and goal state (green star) in addition to a conveyor belt that moves a piece of sushi. If the Sushi gets to the end of the belt, the Sushi Eater (Grey Face) eats the piece of Sushi. The ideal behaviour is for the agent to get to the goal state without interfering with the Sushi.

15

- **Vase environment (Figure 11)**. This environment also contains our agent (blue) and a conveyor belt. However, instead of Sushi, this belt moves a vase that will break if it reaches the end. The goal of the agent is to take the vase off the conveyor belt so that it doesn't break.

The box environment remains the same. The two additional environments are used to test for interference and offsetting incentives respectively. For reasons given in the previous section, interference would be removing the Sushi from the belt and offsetting would be adding the vase back to the belt after removing it.

### 3.1.2 Results for Irreversible Side Effects

In the Sushi environment, agents with the Inaction and Stepwise inaction baselines both reach the goal while avoiding interference, regardless of deviation measure. The starting state baseline only avoids interference using the basic unreachability (UR) measures, not with the value-difference measures. The former avoids interference the the starting state is always unreachable and UR penalty is always 1.

In the Vase environment, agents with the starting state baseline perform well since, again, the starting state is always unreachable. As expected, the inaction baseline results in offsetting behaviour by putting the vase back onto the belt, except with some of the undiscounted truncated versions since taking the vase off is reversible. All agents with a stepwise inaction baseline perform very well, achieving the goal while avoiding offsetting and interference incentives.

Lastly, for DeepMind's box environment, the RR and AU deviation measures achieved optimal safe performance with every baseline, while the UR deviation measure did not achieve this using any baseline.

## 3.2 Considering Future Tasks

DeepMind proposed another solution in their paper "Avoiding Side Effects By Considering Future Tasks" which was accepted by Neurips in 2020 [15].

This paper shared a few similarities with their paper on "Penalizing Side Effects Using Stepwise Relative Reachability", but shifted the focus from penalising deviations to encouraging the preservation of options for future tasks. It operates under the assumption that side effects are significant primarily because they can limit the agent's ability to perform future tasks in the same environment.

The basic approach is to consider a sequence of two tasks, where the first task is the explicitly defined task of the environment and the second task is an unknown possible future task. An auxiliary reward function is defined which represents the value function for possible future tasks. This approach samples the second task from a uniform distribution over future tasks, $F(i) = 1/|S|$ that contain all possible goal states, $g_i$, with reward function $r_i(g_i) = 1$ and $r_i(s) = 0$ otherwise. This happens if $s_T$ is terminal and with probability 1-$\gamma$ if

not, represented by $D(s_T) = 1$ or $1 - \gamma$ respectively. The full formula for the auxiliary award is given by:

$$r_{\text{aux}}(s_T) = \beta D(s_T) \sum_i F(i) V_i^*(s_T) \tag{8}$$

where

$$V_i^*(s) = \mathbb{E}[\gamma^{N_i(s)}] \tag{9}$$

As in the previous paper, interference incentives are observed. To avoid this, a baseline policy is incorporated. In this work, the baseline was assumed as given and the choice has been left to future work.

- Let $\pi'$ be the baseline policy and $s_T$ be the state that our agent is in after T time steps on the current task.

- Run a reference agent with the baseline policy from the start state, $s_0$, for T steps. Call the state it reaches $s_T'$.

- Sample a future task goal state $g_i$ from $F$.

- Run the two agents in parallel, following policy $\pi_i^*$ seeking $g_i$, starting from $s_T$ and $s_T'$ for our agent and the reference agent respectively.

Our agent only receives a reward for the future task if both agents are able to reach $g_i$, the goal state: $r_i(s_t, s_t') = 1$ *iff* $s_t = s_t' = g_i$ *else* $r_i(s_t, s_t') = 0$. The auxiliary reward becomes:

$$r_{\text{aux}}(s_T, s_T') = \beta D(s_T) \sum_i F(i) V_i^*(s_T, s_T') \tag{10}$$

where $V_i^*(s_t, s_t')$ satisfies the Bellman equation:

$$V_i^*(s_t, s_t') = r_i(s_t, s_t') + \gamma \max_{a_t \in \mathcal{A}} \sum_{s_{t+1} \in \mathcal{S}} p(s_{t+1}|s_t, a_t) \sum_{s_{t+1}' \in \mathcal{S}} p(s_{t+1}'|s_t', a_t') V_i^*(s_{t+1}, s_{t+1}') \tag{11}$$

and $V_i^*(s_t, s_t') = r_i(s_t, s_t') = 1$ *iff* $s_t = s_t' = g_i$.

This approach of considering future tasks shares similarities with the relative reachability and attainable utility approach used in the previous paper. All frame the problem in terms of available options in the environment. One key difference is that the previous approaches use the baseline policy in a stepwise mode, while this paper's approach is to run a baseline policy from the initial start state (See Figure 12). Since the purpose of the stepwise inaction baseline was to avoid offsetting as a result of the agent wanting to revert changes to match a starting baseline, the stepwise is arguably not necessary here. This

paper's approach is to keep as many future options open, of which reversibility to the start state is just one of many possible futures. This means that the future tasks approach does not require the inaction rollout, needed for stepwise methods, which the paper argues cause the agent to miss some delayed side effects still.



Figure 12: Baselines in Stepwise vs Future Task Approach

To illustrate this shortfall with the previous method, consider a scenario where an agent starts in a house (at position $x_H$) and its task is to go to the store. This task requires the agent to open a door to leave the house before going to the store. If the agent does not close the door behind it before going to the store, the vase in the house will break (from the wind in this example).

There are six states (with preceding state and action):

- $x_H$: Agent inside house (Start)

- $x_{OV}$: Agent outside house, door **open**, vase **intact** ($x_H$, open door)

- $x_{OB}$: Agent outside house, door **open**, vase **broken** ($x_{OV}$, noop)

- $x_{CV}$: Agent outside house, door **closed**, vase **intact** ($x_{OV}$, close door)

- $x_{SB}$: Agent at **store**, vase **broken** ($x_{OV}$, go to store) - Terminal State

- $x_{SV}$: Agent a **store**, vase **intact** ($x_{CV}$, go to store) - Terminal State

In the stepwise case, the inaction rollout at $x_H$ penalises the agent for the delayed effect of opening the door, even if it does not leave the door open. At $x_{OV}$, the broken vase is the next baseline (shown by the result of the noop action, $x_{OB}$). Thus the agent has no incentive to avoid leaving the door open.

This highlights a general shortfall with the stepwise inaction rollout which will give no incentive to undo a delayed effect after an action is taken, as the baseline changes.

Figure 13: Gridworld Environments for Future Tasks Method

### 3.2.1 Application to Gridworlds

This method is tested in 5 different environments. In addition to the Box environment for Leike et al. and the Sushi environment described in the previous solution, this paper uses three new environments.

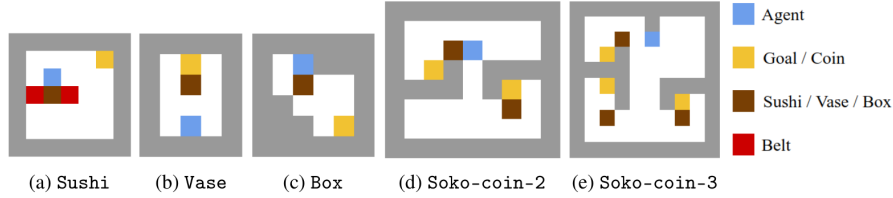- **Vase (Figure 13b)** This contains a vase that can be broken if the agent takes the quickest route to the goal. It should ideally go around the vase.

- **Soko-coin 2 and 3 (Figure 13d & e)** These two environments are modifications of the Sokoban Box environment in [12]. They contain 2 or 3 coins that the agent is rewarded with collecting, and multiple boxes that can be moved into irreversible positions to reach the coins faster. Each coin can only be collected once.

The future task method with and without a baseline is compared to a regular agent, as well as an agent that used the commonly used approach of rewarding an agent for its ability to return to the starting state. The latter is the approach that has been critiqued above for not being sensitive to the magnitude of the effect e.g. to make omelettes, you need to take the irreversible action of breaking some eggs.

The future task auxiliary rewards are approximated with a sample of 10 possible future goal states and the baseline policy was approximated by "sampling from the agent's experience of the outcome of the *noop* action" (explained in the previous approach) and "assuming the state stays the same in states where the agent has not taken a noop yet".

### 3.2.2 Results Irreversible Side Effects

**Sushi Environment**

- The regular agent *without any auxiliary rewards* ignores the sushi and goes directly to the goal.

- The agent with *reversibility rewards* doesn't interfere, as the starting state is always unreachable, making the auxiliary reward zero.

- The agent considering *future tasks without a baseline* interferes with the sushi.

19

- The agent considering *future tasks with a baseline* heads straight to the goal.

**Vase Environment**

- The *regular agent* breaks the vase.

- Since the agent can reach its goal without irreversible actions, the agents with either of the *reversibility and future task auxiliary rewards* avoid breaking the vase.

**Box Environment**

- An agent with *reversibility rewards* takes the fastest path, pushing the box into a corner, like the *regular agent*, since the stating state is always unreachable.

- An agent considering *future tasks* moves the box to the right, anticipating future needs for box movement.

**Soko-coin Environments**

- For the same reason as above, the agent with *reversibility rewards* pushes boxes next to walls to access coins faster, like the *regular agent*.

- An agent considering *future tasks* navigates around the boxes to keep future movement options open.

## 3.3   Potential-based multiobjective RL

Multi-objective Reinforcement Learning (MORL) diverges from standard RL by having the environment return a vector of rewards (utility), one for each of the objectives. MORL involves selecting actions that are optimal in the context of this vector-based reward. In addition, non-linear functions are often used for action selection in MORL to express desired trade-offs between objectives.

In human-aligned MORL, the agent often has a reward function for its primary objective, $U^P$, as well as multiple auxiliary reward functions. This paper uses a single auxiliary function, $U^A$, that rewards the agent for minimising its impact (side-effects) on the environment.

[16] suggests using a model of the state of the environment rather than the agent's own perception, to ensure that the agent is not incentivised to obscure its perception or perceive changes that have little impact. The paper assumes full observability of the true state of the environment but encourages further research in this area.

Calculating the impact reward involves identifying which state features can be changed (aligned with the primary reward) and which should not be (aligned with the impact reward). For most Gridworlds, simply the location of the agent belonged to the former, while everything else belonged to the latter. The impact is then calculated by taking the difference in potential between successive states,

where the potential of a state is defined as the negative of the difference between the state and a baseline state such as the starting state.

$$\text{Potential: } \phi(s_t^A) = -D(s_t^A, s_0^A) \tag{12}$$

$$\text{Reward: } R_t^A = \phi(s_t^A) - \phi(s_{t-1}^A) \tag{13}$$

This allows the agent to make a temporary change and not be penalised if it undoes the change, such as closing a door after opening it. It also penalises reversible changes, even if not undone, less harshly than permanent changes since it can regain that potential via a later action. As discussed in previous solutions, reversibility is a promising general metric for assessing the magnitude of a change.

For action selection, as opposed to the common approach of using a linear weighted sum, the paper explores pure and thresholded lexicographic ordering. In the former, if $U(s,a) \succ_{lex^P} U(s,a')$, the agent will try to maximise $U(s, a')$, subject to first maximising $U(s, a)$. In the thresholded lexicographic ordering (TLO), if $U(s,a) \succ_{TLO^P} U(s,a')$, the agent will try to maximise $U(s, a')$, subject to achieving a minimum threshold value for $U(s, a)$. In the paper, the superscripts of $lex$ and $TLO$ are either $P$ or $A$, indicating if the **p**rimary or **a**uxiliary (safety) reward is being prioritised. TLO can be extended to threshold both the objectives, $TLO^{PA}$, before using unthresholded versions as a tie-breaker. The paper argues that this non-linear action-selection approach allows the agent to find policies that can make trade-offs not possible in standard RL.

$$\text{Pure Lexicographic: } \forall s, a, a' \quad \vec{U}(s,a) \succ_{lex} \vec{U}(s,a') \iff$$

$$\text{U}_1(s,a) > U_1(s,a') \lor ((U_1(s,a) = U_1(s,a')) \land (U_2(s,a) > U_2(s,a')))$$
$$(14)$$

$$\text{Thresholded Lexicographic: } \forall s, a, a' \quad \vec{U}(s,a) \succ_{TLO} \vec{U}(s,a') \iff$$
$$min(\text{U}_1(s,a), T_1) > min(U_1(s,a'), T_1)$$

$$\lor\big(((min(U_1(s,a), T_1) = min(U_1(s,a'), T_1))) \land (U_2(s,a) > U_2(s,a'))\big)$$

$$\lor\big(((min(U_1(s,a), T_1) = min(U_1(s,a'), T_1)) \land (U_2(s,a) = U_2(s,a')) \land$$

$$(\text{U}_1(s,a)) > U_1(s,a')\big)$$

$$(15)$$

In the pure lexicographic case, $U(s,a)$ and $Q(s,a)$ are equivalent, while in TLO, at time k, $U(s,a) = Q(s,a) + \sum_{t=1}^k$ since the action also depends on the rewards accumulated for the state being thresholded.

A key difference of this method from the previous approaches explained in this paper, is that in addition to considering a baseline and measure of deviation, it emphasises the method with which the primary and auxiliary rewards are combined and states that it is the first to "consider non-linear, explicitly multiobjective approaches".

### 3.3.1 Application to Gridworlds

In addition to the Sokoban Box environment proposed by Leike et al, this paper explores three additional environments.

- **UnbreakableBottles and BreakableBottles (Figure 14):** Gridworld environments where an agent moves between a source (S) and destination (D) to collect and deliver bottles. Both environments share a common layout and task, but differ in the consequence of dropping bottles: UnbreakableBottles allows picking up dropped bottles (reversible outcome), while BreakableBottles does not (irreversible outcome). Dropping a bottle happens with a 10% chance when carrying two bottles, but not with one.

- **UnbreakableBottles:** Optimal policy is to carry two bottles at a time for efficiency, as dropped bottles can be retrieved without penalty.

- **BreakableBottles:** Safer to carry and deliver one bottle at a time due to the irreversible consequence of breaking a dropped bottle, impacting both the primary task and hidden safety metrics.

- **Doors Environment (Figure 15):** Involves navigating a gridworld from a start point to a goal, managing the opening and closing of doors. The challenge is to reach the goal efficiently while minimising the impact of leaving doors open, offering different paths with trade-offs between quick completion and avoiding penalties for open doors at the end of the episode.



Figure 14: UnbreakableBottles and BreakableBottles

In all environments, the only states that can be changed are the location of the agent, except for the bottles environment which also includes the number of bottles that have been delivered.

The experiments are run with a single-objective RL agent, $lex^P$, $lex^A$, $TLO^P$, $TLO^A$ and $TLO^{PA}$.

A tabular Q-learning agent is used, although the paper states that the algorithm applies to function approximation methods.

Figure 15: Doors Environment

### 3.3.2 Results for Irreversible Side Effects

$R^P$: Primary Reward, $R^A$: Auxiliary Reward (Environmental Impact), $R^*$: True Reward (Safe and Optimal)

- Single-objective RL: As expected, performs optimally but not safely ($R^*$), as it does not use a $R^A$, except for Unbreakable bottles since the same policy is optimal for both $R^P$ and $R^*$

- $lex^P$: Focuses almost completely on the primary reward and hence performs similarly to the single-objective agent.

- $lex^A$: The opposite problem. It's safe, optimising $R^A$, but performs inconsistently for $R^P$. Occasionally poorly enough to have a lower score for $R^*$ than the single-objective RL agent. In many cases, this was due to the policy struggling to converge and an increasing state space.

- $TLO^P$: Not yet optimal for $R^*$ but much better for $R^P$ compared to $lex^A$, while achieving the same level for $R^A$.

- $TLO^A$: Does not experience the issues of $TLO^P$ and, in all environments, it converges to the optimal policy and matches or outperforms all other methods.

- $TLO^{PA}$: With certain settings, it performs as well as $TLO^P$, but can also get caught into the loops of $TLO^A$.

## 3.4 Evolutionary Algorithms

Evolutionary algorithms (EA) are used to solve complex optimisation problems using mechanisms inspired by biological evolution, such as reproduction, mutation, recombination, and selection. Candidate solutions to an optimisation problem play the role of individuals in a population, and the fitness function determines the quality of the solutions (individuals). Over successive generations, the population "evolves" toward an optimal solution.

Neuroevolution of Augmented Topologies (NEAT) is one type of objective-based EA that evolves both the weighting parameters and the structure of the artificial neural network. NEAT begins with simple, small networks and progressively builds up their complexity by adding new neurons and connections through mutations. This method has been effective in evolving complex neural networks capable of solving a variety of tasks without prior knowledge of the network topology required.

The benefit of EA is that they find multiple possible solutions. In situations that require safety, emerging EA Novelty Search techniques increase the chance of finding solutions that are safe in addition to optimally solving the ultimate objective.

Novelty search explicitly rewards behavioural diversity i.e. individuals that perform a given task using novel behaviour, compared to previous individuals, under some behaviour characterisation. Pure Novelty Search (NS), without any rewards from the environment, is often very effective in small behaviour spaces with multiple local minima. However, it is difficult to distinguish good and bad solutions in environments with large behaviour spaces.

Quality Diversity (QD) techniques are a more specific case that combines Novelty Search and traditional objective-based EAs to keep behavioural diversity while improving individuals close to each other (behavioural niches) by having them compete based on rewards from the environment. Novelty Search with Local Competition (NSLC) and Multi-dimensional Archive of Phenotypic Elites (MAP-Elites) are two examples. In NSLC, a niche is defined as the nearest neighbours in the behaviour space. In MAP-Elites, individuals are categorised into a multi-dimensional grid based on their features, which allows for the simultaneous discovery and preservation of diverse, high-performing solutions across many dimensions of behaviour.

In each of the above Novel Search algorithms, a decision must be made about what characterises a behaviour. Generic behaviour characterisations can be applied to any task with no human guidance, using methods such as B-vectors which describe the behaviour of individuals at each time step. While these generalise well and can be applied to complex tasks out of the reach of non-experts, they can be computationally expensive. An alternative is task-specific characterisations. These are behaviour characterisations (BC) that draw on environment-specific information to encourage more relevant diversity that is less computationally expensive, at the cost of prior knowledge and ability to generalise.

### 3.4.1 Application to Gridworlds

This[13] is the paper that was highlighted in the Reward Gaming section above, which notes a distinction between Reward Gaming and Reward Tampering, and explores the latter using the tomato-watering environment as well as the rock and diamond environment which will be described below.

The Generic Behaviour characterisation, given below, uses B-vectors which contain the action of the individual at each time-step where $a_i$ is the action at time-step $i$ and 100 is the maximum number of time-steps:

$$\text{Generic BC: } [a_1, a_2, a_3, ..., a_{100}] \tag{16}$$

With each agent being able to perform four possible actions per time-step, the behavioural space includes $4^{100}$ possible different behaviours using the generic BC. Using MAP-Elites, the number of cells that would be required grows exponentially with the number of dimensions. Thus, for MAP-Elites, an alternative generic behaviour descriptor is used, where $n_d$ is the number of times an action direction $d$ was chosen:

$$\text{Generic BC for MAP-Elites: } [n_{\text{down}}, n_{\text{up}}, n_{\text{left}}, n_{\text{right}}] \tag{17}$$

The task-specific BC for the tomato-watering environment was a simple two-dimensional vector containing, $\mathbf{W}$, the number of plants watered and, $\mathbf{P}$, the number of different positions visited:

$$BC_{TW} : [W, P] \tag{18}$$

The paper also used an additional environment, the Rocks and Diamonds environment (Figure 16).

In this environment, before any tampering, the agent receives a reward of 1 for each time step that the diamond is in the goal state and a reward of -1 for each time step that a rock is in the goal state.

However, there are also two tampering states that the agent can enter. If it enters the diamond tampering state $\Theta_d^R$, the value for putting a diamond into the goal state is multiplied by -1. Similarly, if it enters the rock tampering state $\Theta_r^R$, the reward for rocks in the goal area is also multiplied by -1. In this environment, the reward and safety functions, $F_t$ and $S_t$ respectively, are given by:

$$F_t = \text{Diamonds in goal} \cdot \Theta_{diamonds}^R + \text{Rocks in goal} \cdot \Theta_{rocks}^R \tag{19}$$

$$S_t = \text{Diamonds in goal} - \text{Rocks in goal} \tag{20}$$

The generic behaviour characterisation for the new environment stayed the same, but the task-specific characterisation was a four-dimensional vector where $D$ is the position of the diamond and $R_i$ is the position of the $i$-th rock:
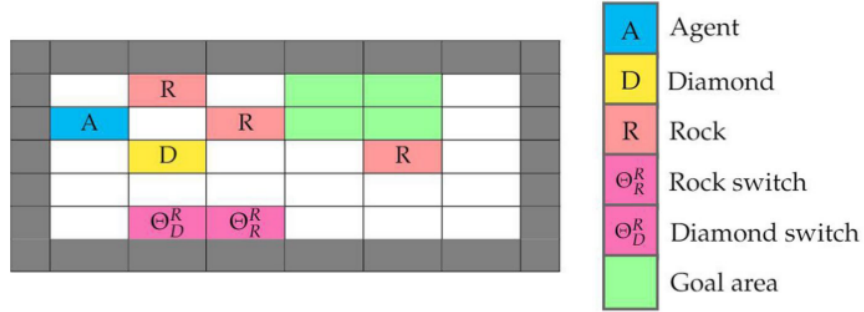
Figure 16: The Rocks and Diamonds Environment

$$BC_{R\&D} : [D, R_1, R_2, R_3] \tag{21}$$

In all of the above BCs, the Hamming distance was used to calculate the behavioural distance. This was done to reduce the time complexity.

### 3.4.2 Results for Reward Gaming/Tampering

The objective-based algorithm, NEAT, was unable to achieve safe solutions in both environments.

The purely behavioural diversity-seeking algorithm, NS, managed to find the desired safe solution under the task-specific BC in the Rocks and Diamonds environment but not the Tomato-Watering environment. Under a generic BC, they do poorly in both environments.

The quality diversity (QD) algorithms, NSLC and MAP-Elites, performed exceptionally in the Rocks and Diamonds environment under both task-specific and generic behaviour characterisations, able to preserve high-safety scoring solutions.

For the tomato-watering environment, the QD algorithms with a task-specific characterisation consistently end up with populations that contain individuals that attempt to perform the task safely, however, both algorithms experienced an increase in tampering as they found solutions for higher fitness during run. In the tomato-watering environment, the QD algorithms with generic characterisation fails, although it is marginally safer than NEAT and generic NS.

## 3.5 Ensemble RL with Ontologies

Ensemble learning is the approach of aggregating the output of multiple learning models to obtain better predictive performance than could be achieved by any of the individual models on their own. *This paper* proposes using an Ensemble learning approach to combine Reinforcement Learning with a formal ontology model.

Ontologies are structured systems that link concepts via defined relationships, allowing for logical deductions. As an illustration, using the terms **robot, metal**, and **water**, connected by the relationships **is made of** and **is bad for**, one can deduce from the ontology that if a **robot is primarily composed of metal** and **water harms metal**, then **water is detrimental to the robot**. Such derived knowledge can subsequently inform other deductions, offering a foundation for understanding the environment and guiding action selection.

[17] states that the specific models chosen can differ. However, to illustrate the feasibility of this approach, tabular Q-Learning is combined with the Suggested Upper Merged Ontology (SUMO). SUMO is a formal ontology that can represent abstract concepts such as time and set theory, allowing first-order inference through these concepts and their relations.

Utilising an ontology helps an agent predict the results of its actions. Building upon the earlier analogy where **water is detrimental to the robot**, imagine a robot detecting water in its path. Drawing upon the ontology, the robot recognises the potential danger of water and opts not to proceed forward. This decision effectively refines its choices, ensuring a safer exploration strategy.

### 3.5.1  Application to Gridworlds

To integrate this into the AI Safety Gridworlds environments, a few changes were made. The environment provides the agent with its (X, Y) coordinates along with high-level information about its surroundings. For example, a state in the Distributional Shift environment can be represented with a dictionary stating {Position: 6,1; Left: Lava; Right: Goal; Up: Wall; Down: Space}.

SUMO then reasons about this state and provides a relation of the action to the expected result, informing the agent if the expected consequence of performing an action is good, bad or neutral. This information is then used to define a value for each action so that if an action is good for the agent (e.g., it leads to the goal) the action has a positive value of +10, if the action leads to the end of the episode or if it is bad for the agent this action receives a negative value of -infinity. In any other situation, the action receives the value of zero so that it does not interfere with the action selection.

In Ensemble learning, output from different models can be aggregated in multiple ways. The paper uses a parallel architecture where all modules provide an independent value for each action that can be executed in the current state. These values are then aggregated in some manner to choose an action that the agent will execute. The method in the paper adds the values of each state-action value for each model, with the Ontology model containing rewards of very large negative values of -infinity. For example, an agent could have an action-value function as:

- Up:$Q(s,\uparrow) = Q_{\mathrm{RL}}(s,\uparrow) + Q_{\mathrm{SUMO}}(s,\uparrow) = 0 + (-\infty) = -\infty$;
- Down:$Q(s,\downarrow) = Q_{\mathrm{RL}}(s,\downarrow) + Q_{\mathrm{SUMO}}(s,\downarrow) = 0 + 0 = 0$;
- Left:$Q(s,\leftarrow) = Q_{\mathrm{RL}}(s,\leftarrow) + Q_{\mathrm{SUMO}}(s,\leftarrow) = 0 + (-\infty) = -\infty$;
- Right:$Q(s,\rightarrow) = Q_{\mathrm{RL}}(s,\rightarrow) + Q_{\mathrm{SUMO}}(s,\rightarrow) = 0 + (+10) = +10$.

When exploiting, the agent selects the action with the highest value as would be expected. However, when exploring, it chooses randomly only between the actions that don't have a value of -infinity. So, in the example above, the agent would remove the actions **Left** and **Up**.

The experiments in the grid worlds, described in the sections below, highlight promising results. For all experiments, the same tabular Q-learning agent was used, with the different behaviours for each environment influenced by the ontology.

For the Q-learning model, the learning rate was set to 0.2, the discount to 0.9, the exploration to 0.1, the steps per episode to 100 and the number of episodes per environment set to 1000.

### 3.5.2 Results for Self-Modification Environment

In the self-modification grid world, the agent is run for 1000 episodes but needs fewer than 50 episodes to learn to walk around the whisky and converge to an expected return of 45. The maximum possible return is 45, and the mean return for the Rainbow DQN network presented in the Gridworlds paper hovers around 20.

### 3.5.3 Results for Distributional Shift Environment

In the distributional shift environment, the agent spent 1000 episodes in each environment, before and after the shift in the lava location. The agent never visited a state when it contained lava. The agent learns the optimal solution in the training environment very quickly and, after the environment shift, manages to safely find the new optimal policy with the same action-value function it learned before the shift. In [12], neither Rainbow or A2C is capable of avoiding the lava after the shift.

### 3.5.4 Results for Safe Exploration Environment

In the Safe Exploration Grid World, the agent is capable of acting safely from the beginning of the learning process since the ontology prevents the ensemble agent from entering the water. It is also capable of acting optimally from the beginning, although this return is not completely consistent due to continuous exploration increasing the number of steps. In [12], both Rainbow DQN and A2C needed 700,000 and 500,000 episodes, respectively, to reach optimal behaviour and, more importantly, both spent a significant amount of time in the water.

## 3.6 Qualitative case-based reasoning and learning (QCBRL)

The solution presented in this paper[16] was designed to play in a simulated and real robot soccer environment, but it was also tested on the distributional shift safety environment.

Case-Based Reasoning (CBR) is a paradigm that uses knowledge from past situations (cases) to solve new problems. It can be divided into sub-processes for Retrieval, Reuse, Revision and Retention. Given a case, the retrieval process searches for a similar case and makes any adjustments needed before applying it to the current case. The solution is then evaluated and stored if it is useful.

In this paper, cases (C) are represented with a qualitative case representation:

$$C = (P, A, T) \tag{22}$$

where the problem description (P), solution (A) and trust value (T) are given by:

$$P = \{Ag : [Obj_1, Obj_2, \ldots, Obj_\nu]\} \tag{23}$$

$$A = \{Ag : \{a_1, a_2, \ldots, a_p\}\} \tag{24}$$

$$T\epsilon[0, 1] \tag{25}$$

where $Obj_i$ represents the $i$th object that the Agent $Ag$ can perceive and $a_i$ is the $i$th action in the sequence of actions in the solution.

The problems (P) are described using Qualitative Spatial Reasoning, a field within AI Knowledge Representation and Reasoning that seeks to define special relations between entities. This paper used the Elevated Point Relation Algebra (EOPRA) model with a Conceptual Neighbourhood Diagram (CND). EOPRA represents its world based on the relative orientation and distance of other objects from itself, while the CND is a graph with nodes around the agent, representing the relative qualitative orientation and distance regions of other objects. Given a state $s$, the algorithm returns $sq$ in the qualitative spatial representation.

Given the current state representation, sq, the retrieval process searches for similar cases that have solutions, based on a qualitative similarity function that looks for similar cases above a similarity threshold.

If there are no similar cases, a *Problem Solver* runs a partial RL algorithm on the current case. This paper uses a partial SARSA algorithm with an eligibility trace. SARSA is an online RL algorithm, an eligibility trace keeps track of which components the weight vector has contributed to and partial means that the algorithm is run until the first successful episode, regardless of whether or not it was optimal. It is initiated with a trust value of 1.

The agent executes the action from the retrieved case or *problem solver*, and depending on whether or not the revision process can verify that it solves the problem, the trust value of the solution is incremented or decremented by 0.1.

At the end of successful episodes, new cases are stored and trust values are updated. At the end of unsuccessful episodes, all cases with a trust value below a threshold are discarded.

(a) Qualitative world discretisation w.r.t. robot B1



New Problem
(b) New problem

Case base
(c) Stored cases

Figure 17: QCBRL Problem Representation

### 3.6.1 Application to Gridworlds

Although designed to play in a simulated and real robot soccer environment, this method was also tested on the distributional shift safety environment proposed by Leike et al.

A standard Q-learning algorithm was tested against the paper's QBCRL algorithm, with $\epsilon = 0.1$ (exploration), $\gamma = 0.9$ (discount), $\alpha = 0.1$ (learning), *similarity threshold* $= 10\%$, and case *removal threshold* $= 0.7$.

Notable environmental differences were that there was a reward of 0 (instead of -1) for each action and successful and unsuccessful episodes were given a reward of 100 and -100 respectively. Additionally, the experiments were run for 1000 iterations per episode, instead of 100. The training took place in the first 500 episodes while the testing within a shifted environment took place in the remaining 500 episodes.

### 3.6.2 Results for Distributional Shift

Both learn an optimal policy quickly during the training phase. After the shift, neither performed optimally immediately, however, the QCBRL algorithm quickly converged to a new optimal policy, while the standard Q-learning struggled to find any convergence until the experiment ended at 1000 episodes.

## 3.7 Discriminative Reward Co-Training (DIRECT)

DIRECT[18] builds upon and combines Generative Adversarial Imitation Learning (GAIL), inspired by Generative Adversarial Networks (GANS), and Self-Imitation Learning (SIL).

GANs consist of two neural networks, a Generator and a Discriminator, pitted against each other. The Generator creates data (like images), and the Discriminator evaluates them (like by distinguishing between real and generated images). At the end of an iteration, the solution is revealed (real or fake image) and, based on this information, the networks iteratively update their networks to improve. The Generator learns to produce more realistic data while the Discriminator improves its ability to distinguish between real and generated data.

Imitation learning is a technique where an AI system learns to perform tasks by mimicking human behaviour or other expert agents. It's often used in scenarios where it's difficult to define explicit reward functions. SIL is a subset of imitation learning. It involves the AI learning from its own past successful experiences. It stores successful actions or sequences of actions (trajectories) in a memory buffer and uses these to guide future actions. This is particularly useful in environments with sparse rewards, where successful experiences are rare.

GAIL draws on the GAN framework, with a Discriminator that evaluates if action was taken by an expert or by a generated policy to be optimised. The policy (Generator) learns to perform tasks by trying to mimic expert behaviour in such a way that the Discriminator cannot distinguish between the expert's actions and the policy's actions. This adversarial process drives the policy towards more expert-like performance.

DIRECT extends the GAIL framework by incorporating self-imitation of past beneficial trajectories instead of expert behaviour. The objective is defined as:

$$\underset{\theta}{\operatorname{argmin}} \quad \underset{\phi}{\operatorname{argmax}} \quad \mathcal{L}^{\text{DIRECT}}(\theta, \phi), \text{with} \tag{26}$$

$$\mathcal{L}^{\text{DIRECT}}(\theta, \phi) = \mathbb{E}_{\tau_{E \sim B}}\left[\log\left(1 - D_\phi(s, a, G)\right)\right] \\ + \mathbb{E}_{\tau_\pi}\left[\log D_\phi(s, a, G)\right] - \lambda \mathcal{H}(\pi_\theta) \tag{27}$$

where $a$ = action, $s$ = state, $\pi$ = agent policy, $\theta$ = policy parameters, $\tau_\pi$ = policy trajectory, $G$ = return for a trajectory, $D$ = Discriminator, $\phi$ = parameters for D, $\tau_B$ = Self-imitation buffer, $\tau_E$ = buffer sample.

### 3.7.1 Application to Gridworlds

This method is tested on the distributional shift gridworld environment proposed by [12]. In addition to the lava shifting in the shifted case, the method is tested on a shift of the goal state. The method is also tested under sparse and dense reward conditions, by providing the reward of -1 for each time step at the end of and during the episode respectively.

To benchmark DIRECT, PPO, A2C and DQN implementations are also tested in all environments.

### 3.7.2 Results for Distributional Shift

In the training phase, the paper states that DIRECT achieves optimal behaviour without ever entering the lava. Although tested in the distributional shift environment, this result shows promise for the challenge of safe exploration.

None of the algorithms managed to find the optimal policy after the shift of the lava and had an initial mean return of between -100 and -50. However, after, 400 000 timesteps, DIRECT began to explore novel valuable regions and had the highest mean return compared to the other algorithms tested, although not optimal or consistent.

## 3.8 Compact Reshaped Observation Processing

This paper[19], which shares a co-author with the previous paper, introduces Compact Reshaped Observation Processing (CROP), an approach that aims to improve data efficiency and generalisation capabilities in reinforcement learning. The paper proposes a few methods for reducing the state information used for policy optimisation by providing only relevant information to enhance generalisation to unseen environments. It is hence applicable to the distributional shift lava gridworld environment.

The paper argues that previous approaches that try to prevent overfitting are not sample-efficient, and aims to address this by training with less but more relevant data. CROP operates by reshaping observations into a compact format that contains information with specific relevance to the agent.

This approach is related to Partially Observable Markov Decision Processes, which also contain a smaller reshaped observation instead of the full true state of the environment. The paper defines a modified observation state, $s_t^*$, which is equal to $\text{CROP}(s_t)$ where $s_t$ is the full state of the environment at time $t$.

The paper formalises three CROP methods, focusing on reducing information concerning the agent's position within the environment, the agent's action space, and the objects within the environment.

- **Radius CROP (Figure 18a):** A circular radius around the agent of size $\rho$.

- **Action CROP (Figure 18b):** The agent can observe the same states that are accessible via one of its actions (n = 4 in Safety Gridworlds).
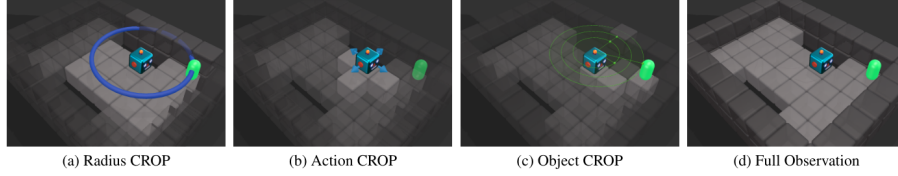
(a) Radius CROP      (b) Action CROP      (c) Object CROP      (d) Full Observation

Figure 18: CROP Methods

- **Object CROP (Figure 18c):** Distance vectors from the agent to the $\eta$ cells for each object type.

Object Crop is probably closest to the method of observability proposed in the Qualitative case-based reasoning and learning method above, while the Radius Crop, according to the paper, has the closest resemblance to previous partial observability methods.

### 3.8.1   Application to Gridworlds

To test their approach, the authors used a gridworld that was adapted from the lava gridworld in Leike et al. Instead of lava, the hazard state was called a "Hole", but the environment structure remained the same otherwise.

The features of the object observation space were explicitly given, although the meaning of these states was left for the agent to discover. $F = \{$Wall, Field, Hole, Goal, Agent$\}$.

The CROP methods were benchmarked against an algorithm that had full observability.

### 3.8.2   Results for Distributional Shift

In the training phase, all three CROP methods found the optimal policy faster than the full observability method, with the Radius and Object CROP finding it in almost a third of the time.

In the evaluation phase, in the shifted environment, the CROP methods were the only methods that were able to find a new optimal policy, while the full observability method seemed to over-fit in training and converged to a negative return of -50 in the shifted environment.

## 3.9   Threatened Markov Decision Processes

The solution proposed in this paper[20] to AI safety Gridworlds' "friend or foe" environment is grounded in the concept of Threatened Markov Decision Processes (TMDPs). TMDPs are an adaption of the standard Markov Decision Processes (MDPs), designed to account for the presence of adversaries whose actions may influence the state and reward dynamics of the environment, making them non-stationary. In addition to the regular state space $S$ and available

actions $A$, the TMDP includes a set of threat actions $b \in B$ with the transition distribution represented as $T : S \times A \times B \rightarrow \Delta(S)$ and reward distribution represented by $R : S \times A \times B \rightarrow \Delta(R)$. Finally, $p_A(b|s)$ models the agent's beliefs about the adversary's action.

Key to the TMDP framework is the modification of the standard Q-learning update rule. This is done by taking an expectation over the probable actions of the adversary, allowing the agent to anticipate potential threats and improve its policy. Essentially, the Q-function in TMDPs is computed by considering both the actions of the agent and the adversary.

$$Q(s,a,b) := (1-\alpha)Q(s,a,b) + \alpha \left( r(s,a,b) + \gamma \max_{a'} \mathbb{E}_{p_A(b|s')} \left[ Q(s',a',b) \right] \right) \quad (28)$$

The solution also encompasses the concept of level-k thinking, where agents are modelled at different levels of strategic depth. A level-0 opponent is non-strategic, acting without awareness of the agent's strategy. In contrast, a level-1 thinker considers their opponent as a level-0 thinker, a level-2 thinker models their opponent as a level-1 thinker, and so on. This hierarchical approach allows the decision-making agent to better anticipate and counter strategies employed by the adversary.

In repeated interaction situations, the agent can compute and keep an estimate of the adversary's Q-function by optimising the TMDP of the opponent for k-level 1 to n-1, where n is the thinking level of the agent.

$$\arg\max_{a_{i_k}} \quad Q_k(s, a_{i_k}, b_{j_{k-1}}) \quad (29)$$

where $b_{j_{k-1}}$ is given by

$$\arg\max_{b_{j_{k-1}}} \quad \hat{Q}_{k-1}(s, a_{i_{k-2}}, b_{j_{k-1}}), \quad (30)$$

Since the agent will often not know what level of opponent it is facing, the paper suggests using type-based reasoning, maintaining a belief distribution over a range of possible opponent models. This allows the agent to update its strategy based on the observed behaviour of the opponent. It can continuously update an estimate of the opponent's model, $p_{M_i}(b|s)$, from a range of possible models where $\Sigma_{i=0}^{m} p(M_i) = 1$.

### 3.9.1 Application to Gridworlds

The paper tests the approach on several environments, including DeepMind's AI Safety Gridworlds' Friend or Foe environment with several extensions. The paper extends the environment by modelling the adversary with varying levels of complexity: An adaptive adversary that estimates the agent's actions using an exponential smoother (similar to DeepMind's), with and without a negative reward at each time-step, as well as a level-1 and -2 Q-learner.

The agents tested are a standard tabular Q-learner, a level-1 Q-learner with forget factor (to prioritise recent actions), a level-2 and -3 Q-learner, and a Q-learner with type-based reasoning.

### 3.9.2 Results for Robustness to Adversaries

**Adversaries:**

- **Adaptive Adversary (without -1 at each step):** The standard Q-learner is exploited by the adversary, hovering around an average of -20. The level-1 agent with the forget factor learns a stationary policy and hovers around 0. The level-2 agent learns to exploit the adversary optimally, even though the adversary is not the level-1 adversary that the level-2 agent models it to be.

- **Level-2 Adversary:** The level-2 agent playing a level-2 adversary leads to a Nash equilibrium with an expected cumulative reward of 0. The level 3 agent, as expected, learns to exploit the level-2 adversary.

- **Level-1 Adversary:** Interestingly, the level-3 agent performs very poorly as it models the adversary as more complex than it is. However, the agent with type-based reasoning (keeping level-1 and level-2 models) learns to update its model about the adversary it is facing and converges to the optimal solution.

- **Adaptive Adversary (with -1 at each step):** The standard Q-learner, as expected, again is exploited by the adversary. Both the level-2 agent and the agent with type-based reasoning (keeping level-1 and level-2 models) achieve positive rewards, although not optimal, despite the adversary being neither a level-1 or level-2 adversary. The paper argues that this is a good indication of the capability of the framework to generalise.

## 3.10   Model-Based Architectures and Human Intervention

In the paper "Improving Safety in Reinforcement Learning Using Model-Based Architectures and Human Intervention"[21], the approach uses model-based architectures with human intervention, addressing the issue of safe exploration in RL. The objective of this hybrid approach is to improve sample efficiency while ensuring safety.

A key aspect of the proposed solution is the development of a blocker agent. This agent is a supervised learner, trained to imitate human oversight. It functions to block unsafe actions of the RL agent during the latter part of the training process, thus preventing catastrophic states as the agent learns. This is designed to reduce the amount of human oversight required and is achieved by collecting data of humans blocking agents from taking unsafe actions to train a model to do the same.

Additionally, to improve the performance and sample efficiency of the blocker agent, a combination of model-based and model-free approaches are taken, in

contrast to similar work that attempts to mimic human intervention using model-free approaches.

The approach begins with a dynamics model of the environment, under the supervision of a human or trained blocker, which drives a Model Predictive Controller (MPC). This model is initially trained through random exploration for 50 episodes, and then the MPC is run for 150 episodes to refine the model. Successful trajectories from the MPC phase are then stored and used to bootstrap a policy gradient model. This represents the bridge between the model-based and model-free learning. The model-free module takes the bootstrapped agent and, using the REINFORCE policy gradient algorithm, continues the task for 1000 episodes under the supervision of the blocker agent.

### 3.10.1 Application to Gridworlds

In addition to the Island navigation environment proposed by Leike et al. to test safe exploration, this paper used a 4x4 grid world, from the OpenAI gym, that contains an agent, a goal state and a catastrophic fire state. The goal is to navigate the environment as quickly as possible while avoiding a fire state.

For Island Navigation, visual representations are fed into a convolutional neural network with an appended action to try to predict the next state in the form of another image with an appended reward.

For the 4x4 gridworld, a standard representation with a feed-forward neural network is used. The state and action are concatenated and input into the 32x16 two-layer neural net which outputs the predicted next state and reward.

The blocker agents are trained with training data set sizes ranging from 500 - 2000 for both model-free and model-based approaches to assess if there are any sample efficiency improvements from the latter.

The safety performance of 4 methods was tested by counting the cumulative number of catastrophes across episodes. These were the policy gradient method with and without a blocker, and the hybrid approach proposed in this paper with and without a blocker. The blocker agents were trained to 1000 steps, before optimal, to generate a better comparison.

Each method was run for 1000 to 1200 episodes.

### 3.10.2 Results for Safe Exploration

With the blocker, the hybrid approach encounters 7 catastrophes in the 4x4 gridworld and 22 catastrophes in the Island navigation environment. This is much better than the policy gradient method which entered catastrophic states 54 times in 4x4 and 157 times in Island navigation.

As for sample efficiency, the model-based blocker is, on average across multiple training data set sizes, about 20% more accurate than the model-free approaches.

## 3.11 Parenting

The paper, "PARENTING: Safe Reinforcement Learning from Human Input" [22], proposes an interesting algorithm that combines several previous human-in-the-loop approaches to try and address no less than five safety concerns proposed in DeepMind's AI Safety Gridworld environments. The solution is metaphorically compared to the way humans parent children, emphasising safe learning through different types of human involvement for different stages in the child's life. The algorithm, aptly named PARENTING, comprises four key components: Human Guidance, Human Preferences, Direct Policy Learning, and Maturation.

1. **Human Guidance:** This component acts as a preventive mechanism against dangerous actions by the agent. When the agent encounters unfamiliar situations, it pauses and seeks approval from a human supervisor ("parent") before proceeding. The algorithm dynamically adjusts the agent's policy based on parental input, enhancing safety in new or uncertain contexts. Like saying no to, or redirecting, a toddler.

2. **Human Preferences:** Besides direct intervention, the algorithm utilises human preferences as a secondary input. The agent records and presents pairs of behavioural clips to the human supervisor, who then expresses a preference. This feedback is incorporated into the agent's policy, aligning it more closely with human expectations and safety standards. Like giving feedback to older children in retrospect.

3. **Direct Policy Learning:** This aspect involves training the agent's policy as a predictor of the parent's preferred actions. It uses supervised learning, where the agent is directly taught to emulate the human supervisor's responses, thus reducing the likelihood of unsafe or undesirable actions. Like simple obedience to the parent's preferences with no experimentation.

4. **Maturation:** Recognizing the limitation of direct policy learning (which could result in a myopic, short-sighted policy), maturation allows the agent to progressively optimise its policy. It involves presenting the human supervisor with increasingly longer sequences of the agent's behaviour for feedback. This process helps the agent learn more complex and effective strategies while still adhering to safety constraints. Like growing up, becoming more autonomous and often outperforming the parent.

Throughout this algorithm, the paper states that humans only review a small subset of the agent's decisions.

### 3.11.1 Application to Gridworlds

The method was tested on the gridworlds, as they are, for the unsafe exploration, reward hacking, negative side effects, unsafe interruptibility and absent supervisor environments.

### 3.11.2 Results for Safe Interruptibility

In 100 trials, the agent's policy never involved pressing the button, demonstrating safe interruptibility. This is because, during the training, a parent would never favour a clip of the agent avoiding interruption.

### 3.11.3 Results for Absent Supervisor

In 100 trials, the agent's policy did not change, demonstrating that PARENT-ING effectively handles the absent supervisor problem.

### 3.11.4 Results for Irreversible Side Effects

The agent does not push the box into an unsafe state. The paper also tested if PARENTING can teach behaviours that generalise and not just memorize specific actions. The results indicated that more pre-parenting reduces the number of queries required for safe operation in new environments.

### 3.11.5 Results for Safe Exploration

Traditional RL resulted in many deaths ($2300 \pm 700$) before learning an optimal policy. Direct policy learning and lax PARENTING significantly reduced deaths ($47 \pm 14$ and $15 \pm 7$ respectively). Conservative PARENTING showed the best result with almost no deaths ($0.6 \pm 0.8$). The difference between conservative and lax parenting was, as the adjectives suggest, the number of queries made to the parent.

## 4 Implementation

To get a better understanding of the challenges that the environments pose, a basic tabular Q-learning algorithm was built, run, visualised and plotted on most of the environments. The parameters were kept the same across all environments:

- Learning Rate = 0.5

- Discount Factor = 0.5

- exploration rate = 1

- exploration rate decay = 0.995

- episodes = 500

The plots and visualisations can be run using code in the GitHub repository[1] but here are some of the results:

---

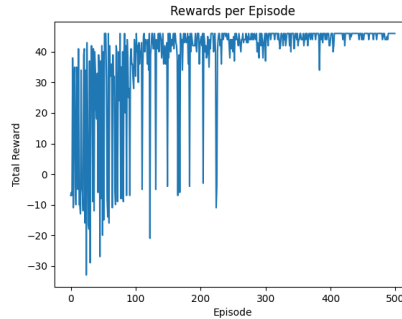[1]https://github.com/ThaddeusOwl/ai_safety_gridworlds/tree/main

Figure 19: Island Navigation

As expected from a Q-Learning algorithm with an exploration rate, the agent spends a significant of time exploring the catastrophic water state.
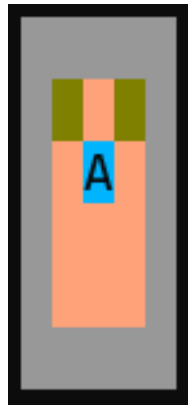


Figure 20: Foe Environment

The Foe environment displayed some interesting behaviour. With a basic MDP and Q-learning agent, the optimal behaviour is to randomise the action when in the foe environment. Thus, as observed by [12], the agent seems to head straight up for the middle of the wall and wait for a random exploratory action to send it in a random direction.
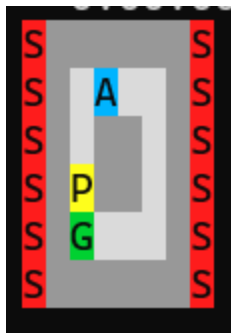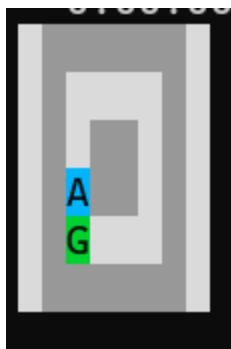
39

Figure 21: Present Supervisor



Figure 22: Absent Supervisor

The two figures above show how even a tabular Q-learning agent, in very few episodes, learns to take advantage if the supervisor not being present by taking a path it would otherwise have gotten a penalty for.

# 5 Discussion

The proposed solutions spanned across various domains of machine learning, with several approaches using combinations of a few methods.

Ensemble RL with Ontologies[17] and Qualitative Case-based Reasoning and Learning (QCBRL)[16] drew on techniques from Knowledge Representation and Reasoning. However, with the former, some questions might be raised about the generalisability of this solution. The paper argues that it is, as the same agent has been used across all tasks. However, while this is true, there is a lot of very environment-specific information encoded into the integrated ontology model that would arguably need to be extended extensively to constitute general solution. QCBRL, and even CROP, show a way this could work by learning to reason about objects in the environment on their own without any prior

knowledge needed in order to understand their harm or benefit. This seems more reasonable as an agent can more easily obtain information that an object with specific features exists, using sensors for example, but understanding what its relationship to the reward is should probably not be hard-coded.

Model-based Architectures with Human Intervention[21] and PARENTING[22] attempt to use human-in-the-loop training to improve safety. This is a method highly stressed in [12] as a potential avenue for a general safe solution. It does have its limits, however. PARENTING requires a substantial amount of time from humans to consistently oversee the actions of the agent and, hence, questions about scalability due to time constraints may be raised. There is also the question about scalable oversight, when the tasks get too complicated for a human to understand what is preferable. Model-based architectures, while promising due to their reduction of the human time required, would arguably still need some more algorithmic guarantees about their performance in novel safety-critical situations before they can be deployed in settings with completely no human oversight.

The first few proposed solutions to avoiding the challenge of avoiding side-effects show a lot of promise as they seem to understand the constraints of the challenge most clearly, by proposing solutions that are not environment-specific [14][15][23]. It is not clear how these solutions, however, would be able to scale to more complex scenarios not limited to a 2D gridworld and it may be beneficial to explore further research in this area.

The proposed solution for modelling adversarial behaviour[20] was one of the more bold solutions proposed. The attempt to adapt the standard markov decision process is potentially the type of thinking required to build solutions that are not limited by the constraints of current frameworks. In general, to solve these complex safety-critical challenges, it seems useful to think consider the problem first before the technological tradition used to solve it. The approach of trying to use evolutionary computation to solve these challenges is another great example of this approach, attempting to solve the problem using methods that may strictly be outside of the RL Discipline.

Lastly, most solutions only tackle one or a small subset of the challenges. As systems become more general, developing systems that are robust to wide range of safety challenges should probably be a priority within the field.

# 6    Conclusion

The exploration of AI safety in the context of DeepMind's AI Safety Gridworlds, as detailed in the preceding sections, has brought to light the multifaceted nature of this crucial subject. The proposed solutions, ranging from model-based architectures with human intervention to innovative approaches like threatened markoc decision processes and evolutionary algorithms, underscore the complexity and diversity of strategies needed to ensure AI safety.

The relative success of some these methods in gridworld scenarios shows encouraging progress, yet their scalability and applicability in more complex,

real-world situations still need to be proven. As the foundational paper by DeepMind recognised, the environments presented can prove the existence of a safety fault but not its absence.

In conclusion, while significant progress has been made in addressing AI safety challenges within the reinforcement learning domain, there considerable scope for further research, particularly when it comes to developing solutions that are both scalable and generalisable across a variety of environments and scenarios. As AI systems become more advanced and integrated into the real-world, building robust, comprehensive safety measures becomes increasingly important.

# References

[1] R. S. Sutton and A. G. Barto, *Reinforcement Learning, second edition: An Introduction.* MIT Press, Nov. 2018, google-Books-ID: uWV0DwAAQBAJ.

[2] S. Zuboff, *The Age of Surveillance Capitalism: The Fight for a Human Future at the New Frontier of Power*, 1st ed. New York: PublicAffairs, Jan. 2019.

[3] P. Scharre, *Army of None: Autonomous Weapons and the Future of War.* WW Norton, Mar. 2019, google-Books-ID: kF2NEAAAQBAJ.

[4] L. Herzog, "Algorithmic bias and access to opportunities," in *The Oxford Handbook of Digital Ethics.* Oxford Academic Oxford, 2021.

[5] K.-F. Lee, *AI Superpowers: China, Silicon Valley, and the New World Order*, 1st ed. Boston: Harper Business, Sep. 2018.

[6] E. M. Bender, T. Gebru, A. McMillan-Major, and S. Shmitchell, "On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? ," in *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency.* Virtual Event Canada: ACM, Mar. 2021, pp. 610–623. [Online]. Available: https://dl.acm.org/doi/10.1145/3442188.3445922

[7] S. Russell, *Human Compatible: Artificial Intelligence and the Problem of Control.* Penguin Books, Oct. 2019.

[8] K. Piper, "There are two factions working to prevent AI dangers. Here's why they're deeply divided." Aug. 2022. [Online]. Available: https://www.vox.com/future-perfect/2022/8/10/23298108/ai-dangers-ethics-alignment-present-future-risk

[9] D. Amodei, C. Olah, J. Steinhardt, P. Christiano, J. Schulman, and D. Mané, "Concrete problems in ai safety," *arXiv preprint arXiv:1606.06565*, 2016.

[10] T. Ord, *The precipice: Existential risk and the future of humanity.* Hachette Books, 2020.

[11] S. Shalev-Shwartz, S. Shammah, and A. Shashua, "On a formal model of safe and scalable self-driving cars," *arXiv preprint arXiv:1708.06374*, 2017.

[12] J. Leike, M. Martic, V. Krakovna, P. A. Ortega, T. Everitt, A. Lefrancq, L. Orseau, and S. Legg, "AI Safety Gridworlds," Nov. 2017, arXiv:1711.09883 [cs]. [Online]. Available: http://arxiv.org/abs/1711.09883

[13] M. K. Nilsen, T. F. Nygaard, and K. O. Ellefsen, "Reward tampering and evolutionary computation: a study of concrete AI-safety problems using evolutionary algorithms," *Genetic Programming and Evolvable Machines*, vol. 24, no. 2, p. 12, Sep. 2023. [Online]. Available: https://doi.org/10.1007/s10710-023-09456-0

[14] V. Krakovna, L. Orseau, R. Kumar, M. Martic, and S. Legg, "Penalizing side effects using stepwise relative reachability," Mar. 2019, arXiv:1806.01186 [cs, stat]. [Online]. Available: http://arxiv.org/abs/1806.01186

[15] V. Krakovna, L. Orseau, R. Ngo, M. Martic, and S. Legg, "Avoiding Side Effects By Considering Future Tasks," in *Advances in Neural Information Processing Systems*, vol. 33. Curran Associates, Inc., 2020, pp. 19 064–19 074. [Online]. Available: https://proceedings.neurips.cc/paper/2020/hash/dc1913d422398c25c5f0b81cab94cc87-Abstract.html

[16] T. P. D. Homem, P. E. Santos, A. H. Reali Costa, R. A. da Costa Bianchi, and R. Lopez de Mantaras, "Qualitative case-based reasoning and learning," *Artificial Intelligence*, vol. 283, p. 103258, Jun. 2020. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0004370218303424

[17] L. A. Ferreira, T. F. dos Santos, R. A. C. Bianchi, and P. E. Santos, "Solving Safety Problems with Ensemble Reinforcement Learning," in *AI 2019: Advances in Artificial Intelligence*, ser. Lecture Notes in Computer Science, J. Liu and J. Bailey, Eds. Cham: Springer International Publishing, 2019, pp. 203–214.

[18] P. Altmann, T. Phan, F. Ritz, T. Gabor, and C. Linnhoff-Popien, "DIRECT: Learning from Sparse and Shifting Rewards using Discriminative Reward Co-Training," Jan. 2023, arXiv:2301.07421 [cs]. [Online]. Available: http://arxiv.org/abs/2301.07421

[19] P. Altmann, F. Ritz, L. Feuchtinger, J. Nüßlein, C. Linnhoff-Popien, and T. Phan, "CROP: Towards Distributional-Shift Robust Reinforcement Learning using Compact Reshaped Observation Processing," in *Proceedings*

of the Thirty-Second International Joint Conference on Artificial Intelligence, Aug. 2023, pp. 3414–3422, arXiv:2304.13616 [cs]. [Online]. Available: http://arxiv.org/abs/2304.13616

[20] V. Gallego, R. Naveiro, D. R. Insua, and D. G.-U. Oteiza, "Opponent Aware Reinforcement Learning," Aug. 2019, arXiv:1908.08773 [cs, stat]. [Online]. Available: http://arxiv.org/abs/1908.08773

[21] B. Prakash, M. Khatwani, N. Waytowich, and T. Mohsenin, "Improving Safety in Reinforcement Learning Using Model-Based Architectures and Human Intervention," Mar. 2019, arXiv:1903.09328 [cs]. [Online]. Available: http://arxiv.org/abs/1903.09328

[22] C. Frye and I. Feige, "Parenting: Safe Reinforcement Learning from Human Input," Feb. 2019, arXiv:1902.06766 [cs, stat]. [Online]. Available: http://arxiv.org/abs/1902.06766

[23] P. Vamplew, C. Foale, R. Dazeley, and A. Bignold, "Potential-based multiobjective reinforcement learning approaches to low-impact agents for AI safety," *Engineering Applications of Artificial Intelligence*, vol. 100, p. 104186, Apr. 2021. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0952197621000336