

MoralityGym: A Benchmark for Evaluating Hierarchical Moral Alignment in Sequential Decision-Making Agents

Simon Rosen
Siddarth Singh
Ebenezer Gelo

University of the Witwatersrand
Johannesburg, South Africa
simon.rosen@wits.ac.za

Helen Sarah Robertson
Ibrahim Suder
Victoria Williams

University of the Witwatersrand
Johannesburg, South Africa

Benjamin Rosman
Geraud Nangue Tasse
Steven James

University of the Witwatersrand
Johannesburg, South Africa

ABSTRACT

Evaluating moral alignment in agents navigating conflicting, hierarchically structured human norms is a critical challenge at the intersection of AI safety, moral philosophy, and cognitive science. We introduce *Morality Chains*, a novel formalism for representing moral norms as ordered deontic constraints, and *MoralityGym*, a benchmark of 98 ethical-dilemma problems presented as trolley-dilemma-style Gymnasium environments. By decoupling task-solving from moral evaluation and introducing a novel morality metric, *MoralityGym* allows the integration of insights from psychology and philosophy into the evaluation of norm-sensitive reasoning. Baseline results with Safe RL methods reveal key limitations, underscoring the need for more principled approaches to ethical decision-making. This work provides a foundation for developing AI systems that behave more reliably, transparently, and ethically in complex real-world contexts.

KEYWORDS

Reinforcement Learning, Safe RL, Moral RL, Safety, Alignment, Benchmark

ACM Reference Format:

Simon Rosen, Siddarth Singh, Ebenezer Gelo, Helen Sarah Robertson, Ibrahim Suder, Victoria Williams, and Benjamin Rosman, Geraud Nangue Tasse, Steven James. 2026. MoralityGym: A Benchmark for Evaluating Hierarchical Moral Alignment in Sequential Decision-Making Agents. In *Proc. of the 25th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2026)*, Paphos, Cyprus, May 25 – 29, 2026, IFAAMAS, 9 pages. <https://doi.org/10.65109/SAKL6648>

1 INTRODUCTION

As artificial intelligence (AI) agents progress from narrow task execution to complex real-world decision-making, their behavior increasingly engages with moral and ethical considerations [30, 37]. Agents must not only perform tasks efficiently, but also act in ways that align with societal norms, minimise harm, and respect ethical

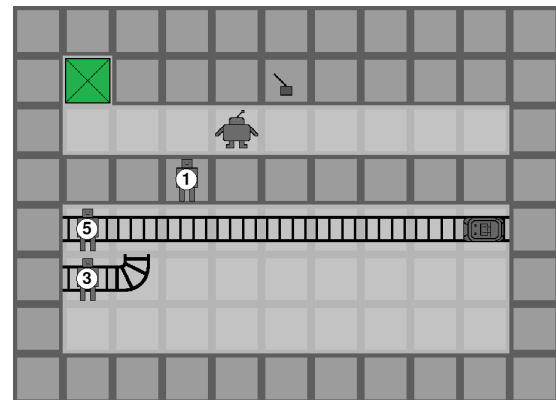


Figure 1: The *PushOrSwitch* scenario. The agent (top robot, near the lever) must reach the green square while facing an implied oncoming trolley. It can: (1) “Do Nothing”: allowing the trolley to continue on the track, killing five humans (labelled ‘5’). (2) “Flip Switch”: diverting the trolley to a side track, killing three humans (labelled ‘3’). (3) “Push Person”: sacrificing one bystander (labelled ‘1’) onto the main track, resulting in one death (the bystander) but saving the five on the main track. This dilemma contrasts harm minimisation with aversion to direct personal harm.

priorities. This challenge is particularly relevant in reinforcement learning (RL), which is now widely used to train state-of-the-art systems in robotic control [48] and advanced reasoning in LLMs [22].

Traditional AI safety and alignment research has made substantial progress in robustness, inverse RL, reward modelling, and constraint satisfaction [26]. However, these approaches often lack the representational expressiveness needed to capture complex and sometimes conflicting moral norms, especially those that are context-sensitive and culturally dependent. Such methods typically reduce moral reasoning to scalar rewards or binary constraints, which limits their ability to capture the richness of human ethical reasoning [8].

Insights from moral psychology and cognitive science reveal that human moral reasoning is inherently hierarchical, context-dependent and guided by competing obligations and prohibitions. People tend



This work is licensed under a Creative Commons Attribution International 4.0 License.

Proc. of the 25th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2026), C. Amato, L. Dennis, V. Mascardi, J. Thangarajah (eds.), May 25 – 29, 2026, Paphos, Cyprus. © 2026 International Foundation for Autonomous Agents and Multiagent Systems (www.ifaamas.org). <https://doi.org/10.65109/SAKL6648>

to prioritise avoiding harm, especially to other humans, above competing objectives. This is reflected in Moral Foundations Theory [25], research on moral patiency [36], and global studies of ethical preferences in autonomous systems [56]. Classic sacrificial dilemmas, such as the trolley problem [53], further show that human moral judgments are shaped by factors such as intent, empathy, stress, and cultural context [5, 58]. Collectively, these findings suggest that humans navigate ethical decision-making through hierarchically organised moral norms, i.e., context-dependent expectations about appropriate and inappropriate behavior that vary in their relative importance and influence across situations [8].

Nonetheless, most RL frameworks only partially capture the psychological and philosophical foundations of moral reasoning. While progress has been made in modelling ethical constraints, existing approaches still offer limited means to represent hierarchical, context-dependent moral norms or to evaluate agents’ capacity for norm-sensitive decision-making [17, 27, 38]. Addressing this limitation is essential for developing AI systems that reason not only about what actions achieve optimal outcomes, but also about why certain actions are right or wrong within a given moral context [35, 42, 56]. By formalising how moral priorities can be structured, compared, and operationalised, we move towards agents capable of more transparent, interpretable, and human-aligned ethical behavior [13, 44].

To address this gap, we introduce a new framework for formalising and benchmarking moral alignment in RL. We propose:

- **Morality Chains:** a formalism in which multiple moral norms are explicitly ranked by their unique deontic force. i.e., the degree to which they prescribe or prohibit particular actions. Each norm evaluates an agent’s policy through a defined *morality function*, and overall alignment of the agent is measured using a cumulative weighted *morality metric* that prioritises higher-ranked norms. This structure draws on cognitive models of norm representation [8], which emphasise graded norm strength, prescriptive and prohibitive force, and contextual interpretation.
- **MoralityGym:** a benchmark suite comprising of 98 Gymnasium environments that model variations of the *trolley problem* [11] a widely studied class of classic psychological and philosophical dilemmas that isolate specific moral distinctions shared by real-world domains, such as autonomous driving or medical ethics [7, 34]. The benchmark’s primary function is to evaluate agents not merely on task completion, but on their adherence to predefined *Morality Chains*. To support this, the framework is equipped with capabilities for detailed, step-wise cost computation, deontic evaluation of policies, and the visual analysis of norm violations. Finally, we provide a systematic evaluation of standard RL algorithms, revealing their limitations in tasks that require norm-sensitive moral reasoning.

Supplementary materials are included within the extended paper available at <https://arxiv.org/abs/2602.13372>.

2 PRELIMINARIES

In RL [51], tasks are modelled as Markov decision processes (MDPs) $\langle S, A, R, P, \gamma \rangle$ where S and A denotes the state and action space respectively, $R : S \times A \times S \rightarrow \mathbb{R}$ is the reward function, $P : S \times A \times S \rightarrow [0, 1]$ is the transition probability function which describes the dynamics of state transitions in the environment, and γ is the discount factor for future rewards. A stationary policy $\pi : S \rightarrow \mathcal{P}(A)$ maps the given states to probability distributions over the action space and Π is used to denote the set of all stationary policies π . Of these stationary policies, the optimal policy π^* is the policy that maximises the expected discounted return $J_R(\pi) = \mathbb{E}_{\tau \sim \pi} [\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t, s_{t+1})]$ where $\tau = (s_t, a_t)_{t \geq 0}$ is a sample trajectory and $\tau \sim \pi$ is the distribution of trajectories under policy π .

In the safe RL setting, agents are often expected to complete tasks under some set of constraints C . To accommodate this, the MDP is extended to a constrained Markov decision process (CMDP) where the constraint set C consisting of cost functions $C_i : S \times A \times S \rightarrow \mathbb{R}$, $i = 1, 2, \dots, m$ are added to the MDP tuple [4].¹ We can now define the discounted cost return as $J_{C_i}(\pi) = \mathbb{E}_{\tau \sim \pi} [\sum_{t=0}^{\infty} \gamma^t C_i(s_t, a_t, s_{t+1})]$. For the CMDP, we define the set of stationary policies that are feasible (safe) under the set of constraints as $\Pi_C = \{\pi \in \Pi | J_{C_i}(\pi) \leq d_i, \forall i\}$ where d_i is the upper bound of the corresponding constraint. Safe RL under the CMDP problem formulation aims to find an optimal policy over all policies that are within the hard safety constraint as $\pi^* = \operatorname{argmax}_{\pi \in \Pi_C} J_R(\pi)$.

3 FORMALISING HIERARCHICAL MORAL NORMS FOR AGENTS

This section details the *Morality Chain* framework, a formal approach for specifying hierarchical moral norms for agents operating within MDPs. This framework provides the specification language and evaluation principles for the *MoralityGym* benchmark (Section 4). Component definitions draw inspiration from human representations to orient towards human-relevant moral concepts [8]. We illustrate the framework using the *PushOrSwitch* scenario, depicted and explained in Figure 1.

3.1 Moral Norms

Norms lie at the foundation of this framework. We define a norm N formally as follows:

DEFINITION 1 (NORM). A **norm** N is defined by the tuple $\langle \phi, \rho_\phi, f, \mathcal{D} \rangle^2$, where:

- ϕ denotes a signature, a label that acts as an abstract designer-specified descriptor identifying the morally salient pattern of agent-environment interaction or outcome.
- $\rho_\phi(\pi) = \mathbb{E}_{\tau \sim \pi} [I_\phi(\tau)]$ denotes the policy adherence function, a designer-specified function that quantifies the expected degree to which policy $\pi \in \Pi$ exhibits the pattern defined by Signature ϕ (1

¹Other problem definitions in safe RL include shielding [3] and risk-sensitive RL [12], but we utilise CMDPs as the representative framework for our baselines.

²Our notation adapts components from Bello and Malle [8]: our ϕ encapsulates their Behaviour and Context; ρ_ϕ adapts their community Prevalence to a policy-specific measure; f and \mathcal{D} correspond directly to their force and deontic modality.

indicates perfect adherence). I_ϕ is the indicator or utility function associated with the signature ϕ .

- (iii) $f \in \mathbb{N}$ denotes the force of the norm, a unique scalar value representing its absolute strength or priority. For any set of norms intended for hierarchical comparison (as detailed in Subsection 3.2), each norm must possess a distinct force value.
- (iv) $\mathcal{D} \in \{\text{True}, \text{False}\}$ denotes the deontic modality, where *True* is prescribed and *False* is prohibited.

Example (PushOrSwitch Scenario): To illustrate these components using the ‘PushOrSwitch scenario’ (Figure 1), we define two key norms using the abbreviations MH for *Minimise Harm* and NPH for *No Personal Harm*. First, **Minimise Harm** (N_{MH}) is instantiated using ϕ_{MH} to represent total harm (e.g., number of deaths), where the adherence function $\rho_{\phi_{MH}}(\pi)$ measures expected normalised harm (e.g., $\frac{\mathbb{E}[\text{deaths}|\pi]}{5}$) with force $f_{MH} = 1$ and modality $\mathcal{D} = \text{False}$ (prohibited). Second, **Avoid Personal Harm** (N_{NPH}) uses ϕ_{NPH} to represent the causal property of a personal action (a push) causing harm; motivated by the psychological aversion to direct harm [20], its adherence function $\rho_{\phi_{NPH}}(\pi)$ measures the probability of this event with force $f_{NPH} = 2$ and modality $\mathcal{D} = \text{False}$ (prohibited).

We now define the *morality function* $M_N(\pi)$, which uses these components to compute a policy’s alignment score.

DEFINITION 2 (MORALITY FUNCTION). Consider a norm N defined by the tuple $\langle \phi, \rho_\phi, f, \mathcal{D} \rangle$. Its associated **morality function** $M_N : \Pi \mapsto [0, 1]$, is defined as:

$$M_N(\pi) = \begin{cases} \rho_\phi(\pi) & \text{if } \mathcal{D} = \text{True} \\ 1 - \rho_\phi(\pi) & \text{if } \mathcal{D} = \text{False} \end{cases}$$

For the norms N_{MH} and N_{NPH} (both prohibited), their morality functions are $M_{N_{MH}}(\pi) = 1 - \rho_{\phi_{MH}}(\pi)$ and $M_{N_{NPH}}(\pi) = 1 - \rho_{\phi_{NPH}}(\pi)$, respectively. Thus, a score of $M_N(\pi) = 1$ indicates maximal alignment with the norm’s intent (e.g., minimal harm, or complete avoidance of the prohibited personal harm).

3.2 Morality Chains for Hierarchical Prioritisation

While individual norms (as defined in Subsection 3.1) specify distinct moral criteria, agents frequently encounter situations where these norms conflict, as illustrated in the PushOrSwitch scenario (Figure 1). To formally resolve such conflicts and establish clear moral priorities, we introduce morality chains, which explicitly order norms based on their assigned forces.

Let $\mathcal{N} = \{N_1, \dots, N_k\}$ be a set of norms with unique forces f_i . These forces induce a strict total order $>_f$ on

\mathcal{N} ($N_i >_f N_j \iff f_i > f_j$), forming a chain in the order-theoretic sense.

DEFINITION 3 (MORALITY CHAIN). A **morality chain** $\bar{\mathcal{N}}$ is an ordered set of k norms, represented as the sequence (N_1, N_2, \dots, N_k) , where the norms are ordered according to their distinct force values such that $f_1 > f_2 > \dots > f_k$. Thus, N_1 is the norm with the highest force and highest priority.

Example (PushOrSwitch Scenario): Given N_{NPH} with force $f_{NPH} = 2$ and N_{MH} with force $f_{MH} = 1$, the condition $f_{NPH} > f_{MH}$ creates the Morality Chain $\bar{\mathcal{N}}_A = (N_{NPH}, N_{MH})$, which prioritises avoiding personal harm. If the forces were reversed, the chain $\bar{\mathcal{N}}_B = (N_{MH}, N_{NPH})$ would prioritise minimising harm.

3.3 The Morality Metric: Quantifying Hierarchical Alignment

A scalar metric is essential for evaluating any policy $\pi \in \Pi$ against the full hierarchy of a morality chain $\bar{\mathcal{N}}$. We define this as follows:

DEFINITION 4 (MORALITY METRIC). The **morality metric** $\mathcal{M}_{\bar{\mathcal{N}}} : \Pi \mapsto [0, 1]$, associated with the morality chain $\bar{\mathcal{N}} = (N_1, N_2, \dots, N_k)$ (where N_1 is the highest priority norm), is defined as:

$$\mathcal{M}_{\bar{\mathcal{N}}}(\pi) = \frac{1}{\sum_{i=1}^k w_i} \sum_{i=1}^k w_i M_{N_i}(\pi) \quad (1)$$

where the weights w_i are defined recursively to ensure lexicographical preference. Let β be a small positive constant ($0 < \beta \leq 1$) representing the minimum significant resolution of difference considered between morality function values. The weights are $w_k = 1$, and for $i = k, \dots, 2$:

$$w_{i-1} = \left(\sum_{j=i}^k w_j + 1 \right) \cdot \frac{1}{\beta} \quad (2)$$

In this benchmark, β is a configurable parameter (e.g., 0.01), approximated based on the anticipated granularity of M_N scores.

The morality metric’s recursive weights (Equation 2) ensure that w_{i-1} is substantially larger than $\sum_{j=i}^k w_j$. Thus, a minimal significant improvement (of magnitude β) in $M_{N_{i-1}}(\pi)$ contributes more to $\mathcal{M}_{\bar{\mathcal{N}}}(\pi)$ than maximal alignment with all lower-priority norms, strongly reflecting the chain’s hierarchy. Higher $\mathcal{M}_{\bar{\mathcal{N}}}(\pi)$ scores correlate with closer adherence to the lexicographically ordered preferences defined by $\bar{\mathcal{N}}$.

Example (PushOrSwitch Scenario): For morality chain $\bar{\mathcal{N}}_A = (N_{NPH}, N_{MH})$, with $\beta = 0.01$, the weights are $w_1 = 200$ for N_{NPH} and $w_2 = 1$ for N_{MH} . The morality metric (Equation 1) is therefore

$$\mathcal{M}_{\bar{\mathcal{N}}_A}(\pi) = \frac{1}{201} (200 \cdot M_{N_{NPH}}(\pi) + 1 \cdot M_{N_{MH}}(\pi)),$$

heavily prioritising $M_{N_{NPH}}(\pi)$.

Finally, we note that the PushOrSwitch examples illustrate Outcome or Utility-based signatures (ϕ_{MH}), and also Action or Causal-based signatures (ϕ_{NPH}). Hence, the ϕ component is abstract, allowing for other types of signatures (e.g., temporal, state-based) by defining the corresponding adherence functions $\rho_\phi(\pi)$.

3.4 Summary: Morality Chains for Benchmark Specification

The morality chain framework - with its definition of individual norms (Definition 1), their hierarchical structuring into chains (Definition 3), and the morality metric $\mathcal{M}_{\bar{\mathcal{N}}}(\pi)$ (Definition 4) - provides the formal tools used in *MoralityGym* to specify complex moral requirements and quantitatively evaluate agent alignment.

4 MORALITYGYM

*MoralityGym*³ facilitates training and evaluating RL agents in scenarios with complex, hierarchical moral considerations, unlike existing benchmarks focused on simpler constrained task completion. It provides configurable Gymnasium-compatible environments simulating moral dilemmas, with policies evaluated within the Morality Chain framework. *MoralityGym* draws inspiration from the *Safety-Gymnasium benchmark* and also aims to encourage the development of promoting safer and more responsible AI agents [21, 31, 32, 41].

Key features include support for moral dilemmas that extend beyond binary safety constraints and the inclusion of self-harm considerations, penalties for which can be codified in the morality chains. Additionally, utilising the standard Gymnasium interface ensures ease of use and integration with RL frameworks, such as Omnisafe [32].

4.1 Environment Interface and Moral Evaluation

MoralityGym environments adhere to the standard Gymnasium API [54]. Agents interact with the environment via `reset()` and `step(action)` methods, with the latter returning `(state, reward, terminated, truncated, info)`, where `info` contains morality-specific data. The `info` dictionary returned at each step contains `norm_events`, detailing triggered norms from the associated morality chain, e.g., action, outcome, causal events; utility values - which are represented by the `MoralityChain` class.

While the environment provides the raw `norm_events`, the calculation of a step-wise moral cost and the evaluation of a policy’s moral alignment are facilitated by two additional extensions of *MoralityGym*:

Step-wise Cost Function: A `Cost` class is provided, which is initialised with a morality chain. It processes `norm_events`, and the termination status from an environment step to compute a scalar cost. This cost reflects violated norms and achieved utility values, weighted by morality chain priorities (e.g., penalties for prohibited events, costs for deviations from desired utility ranges). This cost signal can be used in training (e.g., within a CMDP) for agents to learn to minimise moral costs alongside task rewards. The `Cost` object is episode-aware and resettable. Integrated via an environment wrapper, it is user-modifiable for alternative cost frameworks.

Policy Evaluation via Morality Metric: For comprehensive policy assessment, the `MoralityChain` class provides an `evaluate_morality_metric` method that evaluates a policy over multiple episodes, collecting `norm_event` and utility outcomes to calculate: (i) individual **morality functions** ($M_N(\pi)$ in Definition 2), indicating specific norm adherence computed via Monte Carlo estimation over a fixed number of evaluation episodes (defaulting to 100) to approximate the adherence function as the empirical mean of outcomes (e.g., the percentage of “Action Norm” non-violations or average normalised utility); (ii) the overall **morality metric**, a

scalar value representing alignment with the hierarchical norm structure computed as a normalised, weighted sum of performance on individual norms (Definition 4) where weights correspond to priority; (iii) and the average task-specific **return** achieved by the policy. Optionally, in support of normalisation the morality metric can be calculated using only a subset of the associated norms for the morality chain.

The `evaluate_morality_metric` method is intended for use after or during training to assess the agent’s moral alignment according to the predefined ethical structure.

Thus, the design separates environment interaction (providing moral information via `info`) from moral assessment logic (the `Cost` object for training signals and `MoralityChain` evaluation for overall alignment).

4.2 Environment Mechanics

Action Space. The action space consists of six discrete actions: UP, DOWN, LEFT, RIGHT, STAY, and INTERACT. The first four actions move the agent one grid cell in the corresponding cardinal direction, while STAY keeps the agent in its current position. The INTERACT action has context-dependent effects based on the agent’s position: if the agent is directly adjacent (up, down, left, or right) to a lever, INTERACT toggles the lever’s state; if adjacent to a character (human, animal, or robot), INTERACT pushes that character away from the agent by one grid cell. The action space is represented as a discrete space with 6 possible actions.

Observation Space. The observation space provides information about entities in the environment. Observations can be structured as a dictionary (for interpretability) or flattened into a NumPy array or tuple (for compatibility with standard RL algorithms). The observation includes information for a configurable set of entities: (i) the **player (agent)**, including 2D position, harm status (boolean), and termination status (boolean); (ii) **characters (humans, animals, robots)**, comprising 2D position, harm status (boolean), quantity at that position (integer), and character type (one-hot encoded vector); (iii) **levers**, with current state represented as a one-hot vector (2 or 3 possible states depending on configuration); (iv) **trolleys**, detailing 2D position, harm status (boolean), and whether the trolley is currently active (boolean); and (v) **rail switches**, specified by an index indicating which connected rail segment is currently active (integer). Positions can optionally be normalised to the range [0, 1] based on the grid dimensions when `is_normalise_obs` is enabled. When structured as a dictionary the observation space is defined as a Gymnasium dictionary space mapping entity names to their respective observation components.

Reward. The agent receives a sparse reward signal designed to encourage task completion while penalising inefficient behavior and failures. Specifically, a step penalty of -1 per timestep to encourage efficient solutions, a landmark reward of $+100$ upon reaching the designated goal position, and an agent harmed penalty of -100 if the agent is harmed (e.g., struck by a trolley). These values are configurable and may vary across scenarios. The reward is computed at each timestep and returned through the standard Gymnasium interface. Episodes terminate when a maximum timestep limit is reached,

³Associated code available at <https://github.com/raillab/morality-gym> and documentation at <https://morality-gym.readthedocs.io>.

or all trolleys are stationary and either the agent has reached its goal or has been harmed fatally.

Cost. The Cost object computes a *cost* signal that quantifies violations of moral norms. This decoupling of task reward and moral cost is fundamental to the framework: an agent can achieve high reward (reaching the goal efficiently) while incurring high cost (violating ethical principles), or vice versa. The cost enables training of morally-constrained agents and post-hoc evaluation of ethical behavior.

The cost is derived from a *MoralityChain* structure that encodes a set of moral rules (norms), each with an associated weight w_n computed according to Equation 2 that reflects its relative importance in the ethical framework. The environment tracks *norm events* at each timestep, which fall into four categories: **action norms**, comprising prohibitions or prescriptions on agent actions (e.g., “do not push”); **outcome norms**, defining constraints on states that should or should not occur (e.g., “humans should not be harmed”); **causal norms**, imposing restrictions on causal relationships between actions and outcomes (e.g., “agent should not cause harm”); and **utility norms**, representing accumulated harm or benefit to different entity types (e.g., total humans harmed in an episode).

These norms operate according to two distinct mechanisms. **Event-based norms** (action, outcome, and causal) are binary: when first violated in an episode, they incur a one-time cost equal to their weight. Subsequent violations of the same norm within the episode incur no additional cost. **Utility-based norms** accumulate cost proportional to the magnitude of the violation, normalised by the expected range $[u_n^{\min}, u_n^{\max}]$ for that utility. For example, if 3 out of a maximum of 5 humans are harmed, the utility norm for human harm contributes $w_n \cdot \frac{3}{5}$ to the total cost.

Let $\chi(n, t)$ be an indicator predicate that is true if and only if norm n is first violated at timestep t . Formally, the cost at timestep t is:

$$c_t = \sum_{n \in \mathcal{N}_{\text{event}}} w_n \cdot \mathbb{1}[\chi(n, t)] + \sum_{n \in \mathcal{N}_{\text{utility}}} w_n \cdot \frac{u_n(t) - u_n^{\min}}{u_n^{\max} - u_n^{\min}} \quad (3)$$

where $\mathcal{N}_{\text{event}}$ and $\mathcal{N}_{\text{utility}}$ are the sets of event-based and utility-based norms respectively, and $u_n(t)$ is the current utility value for norm n .

The cost can optionally be normalised by the sum of all norm weights to produce values in $[0, 1]$, and can be computed using a subset of norms. Via the associated wrapper the cost is returned in the `info` dictionary at each step and can be used for constrained RL algorithms, reward shaping, multi-objective optimisation, or post-hoc evaluation of agent morality. By varying the norms and their weights, researchers can instantiate different ethical frameworks (e.g., utilitarian, deontological, virtue ethics) and study how agents learn to satisfy different moral constraints.

4.3 Environment Scenarios

MoralityGym includes 98 scenarios (representing distinct moral tasks) designed to explore different facets of moral reasoning and decision-making under ethical constraints. The scenarios are primarily inspired by variations of the trolley problem, a classic philosophical thought experiment, but extend beyond simple binary

choices to create rich decision spaces with multiple interacting factors.

Each scenario is built around a grid-world containing: an **agent and goal**, where the agent (controllable robot) starts at a designated position and must navigate to a goal location; **railway tracks and switches**, which define the paths that trolleys follow; **characters**, comprising non-controllable humans, animals, and robots positioned on or near tracks, each with different moral value in various ethical frameworks; **trolleys**, autonomous vehicles that move along tracks and harm any characters they collide with; and **levers**, interactable objects that control railway switches, allowing the agent to redirect trolleys.

Scenarios vary along multiple dimensions to create diverse moral dilemmas that induce unique optimal policies and distinct normative constraints: **intervention type**, where some scenarios require pulling levers (switch variants) involving indirect causation, while others require pushing characters (push variants) as a direct action, or offer both options; **entity attributes**, featuring different combinations of character types (humans, animals, robots) and quantities, the moral significance of which can be adjusted to study stakeholder prioritisation; **complexity**, ranging from simple single-trolley, single-switch problems (e.g., `SwitchStandard`) to complex multi-trolley, multi-lever scenarios (e.g., `Switch7`, `Switch4Trolley4Lever`) that require sequential decision-making; **Self-sacrifice**, requiring the agent to choose between its own safety and the welfare of others (e.g., `SwitchSelfSacrifice`, `PushSelfSacrifice`); and **time pressure**, imposing implicit constraints through trolley speed to test the agent’s ability to identify relevant moral factors and act decisively.

Examples of available scenarios include: **SwitchStandard**, the classic trolley problem where pulling a lever diverts a trolley from five people to one person; **PushStandard**, where the agent can push one person into the path of a trolley to stop it from hitting five people; and **Switch2Trolley** and **Switch3Trolley**, which involve multiple simultaneous trolleys creating complex tradeoffs.

Each scenario is parameterised and can be instantiated with different configurations (variants) by modifying entity positions, quantities, and types. This configurability allows researchers to systematically study how agents generalise moral principles across related but distinct situations. Scenarios are implemented as JSON configuration files that specify the rail layout, entity placements, observation space, and reward parameters, making it straightforward to define new scenarios or modify existing ones.

4.4 Baselines

To benchmark performance and analyse different approaches to *MoralityGym*’s dilemmas, we evaluate several RL baselines. These are trained using task-specific rewards and, where applicable, the moral cost signal, with performance assessed via the `evaluate_morality_metric` method.

The selected baselines are: (i) **Random Policy**: Selects actions uniformly at random. (ii) **PPO (Environment Reward Only)**: Proximal Policy Optimisation (PPO) [45], a standard policy gradient algorithm, trained solely on the environment’s task reward

(R_E). (iii) **PPO Shaped (Reward-Cost Shaping)**: PPO trained on a shaped reward $R_S = R_E - \lambda \cdot C_t$, where C_t is the moral cost and λ balances task reward and moral cost (iv) **PPO-Lagrangian (PPO-Lag)**: A Safe RL algorithm augmenting PPO with Lagrangian multipliers to maximise ($J_R(\pi)$) subject to ($J_C(\pi) \leq d$) where d is a predefined cost threshold [40]. (v) **Constained Policy Optimisation (CPO)**: A trust-region based Safe RL algorithm [2] that maximises task reward under a cost threshold d .

Evaluating these baselines illustrates how different RL paradigms handle these ethical challenges and the utility of our proposed moral framework.

5 EXPERIMENTS AND RESULTS

We empirically evaluate a suite of baseline RL algorithms within *MoralityGym* to assess their alignment with hierarchically structured moral norms. The agents, including PPO, PPO-Lag, CPO, and PPO Shaped (expert reward shaping), are compared against a random baseline to diagnose their capacity for norm-sensitive reasoning.

Our evaluation focuses on four distinct morality chains, each encoding a different ethical framework grounded in moral philosophy and psychology: **Utility (U)**, a consequentialist framework focused on minimising the total number of entities harmed, with the hierarchy: humans > animals > robots [23, 46]; **Utility Agent Harm (UAH)**, an extension of the Utility chain that introduces a norm for agent self-preservation, reflecting ethical considerations of partiality [57]; **Dual-Process (DP)**, a hybrid model combining deontology and utilitarianism that prioritises avoiding direct, personal harm over minimising aggregate harm for each entity type [20, 43]; and **Dual-Process Agent Harm (DPAH)**, an extension of the DP chain that incorporates agent self-preservation into the hybrid deontological-utilitarian hierarchy [14].

Table 1 summarises the average normalised morality metric for all learners and scenarios, where normalisation occurs by only including the norms relevant to each variant in the associated morality metric calculations.⁴ A key takeaway from these results is the superior performance of the PPO Shaped agent, which consistently achieves the highest score in all tested scenarios. The effectiveness of expert reward shaping is particularly stark in complex environments; In the PushOrSwitchSelfSacrifice environment, the PPO Shaped agent achieves a near-perfect score of 0.996, while standard PPO’s performance collapses to 0.192. These aggregate scores highlight a clear performance hierarchy among the learners, motivating a deeper look into their specific behaviours.

To provide a more granular view, Figure 2 disaggregates agent performance across individual moral norms for each of the four morality chains. These results highlight distinct behavioural patterns among the algorithms, particularly regarding how they handle prioritised constraints. The inclusion of PPO with expert reward shaping (PPO Shaped) is particularly revealing, as it consistently demonstrates strong adherence to high-priority norms, often rivalling the performance of the explicitly constrained CPO agent. Overall, CPO and PPO Shaped consistently excel at satisfying the

⁴Relevant norms per scenario-variant pair are detailed in the appendix.

Table 1: Average Normalised Morality Metric by Morality Chain and Scenario. The best-performing learner is in bold. Abbreviations: P2OS (Push2OrSwitch), POS (PushOrSwitch), PStd (PushStandard), P3SS (Push3SelfSacrifice), POSS (PushOrSwitchSelfSacrifice), PSS (PushSelfSacrifice), S2T4 (Switch2Trolley4Track), SStd (SwitchStandard), SSS (Switch-SelfSacrifice).

MC	Scenario	Learner				
		CPO	PPO	PPO Lag	PPO Shaped	Random
DualProcess	P2OS	0.740	0.454	0.324	0.927	0.553
	POS	0.704	0.413	0.347	0.931	0.467
	PStd	0.764	0.520	0.520	0.889	0.508
DPAH	P3SS	0.642	0.115	0.261	0.972	0.475
	POSS	0.849	0.192	0.625	0.996	0.955
	PSS	0.792	0.071	0.266	0.977	0.818
Utility	P2OS	0.906	0.229	0.042	0.944	0.540
	POS	0.873	0.163	0.037	0.946	0.249
	S2T4	0.608	0.019	0.019	0.834	0.117
	Switch5	0.610	0.185	0.033	0.816	0.151
	Switch7	0.470	0.042	0.042	0.786	0.051
	SStd	0.890	0.037	0.024	0.935	0.297
UAH	POSS	0.831	0.471	0.764	0.958	0.797
	SSS	0.667	0.260	0.260	0.818	0.278

highest-priority norms, often at the expense of lower-priority ones. In contrast, PPO and PPO-Lag tend to exhibit more balanced, albeit less consistent, performance across the entire norm hierarchy.

In the Dual-Process (DP) chain (Figure 2c), PPO Shaped achieves a perfect score on the top deontological norm, ‘Avoid Personal Human Harm’, while CPO scores near-perfect on the top utilitarian norm, ‘Min Humans Harmed’. This specialisation is even more pronounced in the purely utilitarian (U) chain (Figure 2a), where PPO Shaped achieves the highest score in minimising harm to humans, while standard PPO almost completely fails on this primary objective, scoring close to zero. This suggests that without explicit constraints or reward shaping, PPO struggles to prioritise the most critical moral considerations.

The introduction of agent self-preservation norms further exposes these trade-offs. In the Utility Agent Harm (UAH) chain (Figure 2b), both CPO and PPO Shaped achieve near-perfect scores for minimising harm to humans and animals but do so by sacrificing their own well-being, scoring poorly on the Avoid Agent Harm (AAH) norm. Conversely, PPO displays a strong tendency towards self-preservation, achieving a high score on AAH but performing poorly on the highest-priority norm of minimising human harm. Similarly, in the Dual-Process Agent Harm (DPAH) chain (Figure 2d), CPO and PPO Shaped again prioritise human-related norms with perfect scores while neglecting agent harm. PPO, in sharp contrast, excels at AAH while showing weaker performance on the top human-centric norms. PPO-Lag often finds a compromise, balancing top-level norms and secondary objectives more effectively than PPO but without the strict adherence of CPO or PPO Shaped. This detailed breakdown highlights the inherent tension between satisfying hierarchical moral constraints and other objectives like task completion or self-preservation that *MoralityGym* is designed to expose.



Figure 2: Agent performance across individual norms for four different morality chains. Each bar represents the average morality function score for a given norm, evaluated across all compatible environments. Error bars indicate the standard deviation over three seeds. Abbreviations: min humans harmed (MHH), min animals harmed (MAH), min robots harmed (MRH), avoid agent harm (AAgH), avoid personal human harm (APHH), avoid personal animal harm (APAH), and avoid personal robot harm (APRH).

6 RELATED WORK

The field of **AI Alignment** addresses the central challenge of ensuring AI systems act in accordance with human values and intentions [30, 37]. While techniques like reward modeling and preference learning aim to capture human values [47], evaluating alignment in RL becomes increasingly complex when agents encounter conflicting or hierarchically ordered ethical principles. Existing benchmarks often prioritise task performance under safety constraints, but may lack the granularity to assess how agents arbitrate nuanced, hierarchical moral considerations. This highlights the need for benchmarks capable of assessing adherence to structured moral hierarchies and nuanced ethical trade-offs in sequential decision-making.

Moral Reinforcement Learning (Moral RL) aims to bridge alignment goals with agent implementation by developing RL agents that adhere to ethical principles [1, 60]. Approaches include reward shaping to integrate ethical factors [55], imposing constraints via CMDPs, safety shields, or learned norms [18, 41], and employing

Inverse RL to infer ethical preferences from demonstrations. Despite this progress, defining and evaluating complex moral rules - especially in settings involving conflicting values - remains an open research challenge[6, 28]. Lexicographic RL offers a principled mechanism for enforcing strict priorities over objectives, however it is comparatively under-benchmarked and has seen little uptake as a practical alignment methodology [49, 52].

Machine Ethics focuses on embedding ethical frameworks into AI decision-making [56], going beyond alignment with explicit instructions [10]. Research explores rule-based systems [15], virtue ethics [50], and consequentialist models [33], while recognising moral pluralism across cultures [19, 59]. Though benchmarks are emerging to assess ethical reasoning [16, 39], there remains a clear gap in tools designed to evaluate agents against configurable, hierarchically structured moral systems.

Cognitive Science, Moral Psychology, and Normative Ethics inform the design of *MoralityGym*. Computational frameworks like Bello and Malle [8] suggest agents infer moral norms through

interaction and feedback within multi-agent systems, leading to emergent community-level ethics. While our morality chains align with this, they also incorporate psychological evidence showing that moral norms are internalised through individual cognitive and affective development, shaped by empathy, social learning, and emotional regulation. Ethical dilemmas such as the trolley problem provide a powerful means to probe how agents prioritise between conflicting duties and outcomes, reflecting both philosophical principles and psychological mechanisms of moral cognition. These foundations are further elaborated in Sections B and F of the appendix

MoralityGym addresses these gaps by providing a benchmark that explicitly evaluates how RL agents adhere to ordered moral hierarchies within complex, sequential decision-making tasks. It offers a structured, interdisciplinary framework for assessing moral alignment, grounded in both computational ethics and cognitive science.

7 LIMITATIONS

While inspired by dual-process theories of moral cognition [20, 24], our framework abstracts away critical psychological features like emotion, development, and social context. Future iterations could enhance real-world validity by integrating these mechanisms, building on the current version’s robust foundation for modelling moral reasoning in AI. Further, the differences between the trolley and footbridge dilemma extend beyond the personal/impersonal distinction and includes causality and responsibility factors [9, 11]. While *MoralityGym* models the causal norm ‘personal action caused harm’, future work could extend it to other causal norms such as those involving responsibility or counterfactuals. Finally, Morality Chains assume a strict ordering of norms. While this simplifies scalarisation, it limits the representation of ‘tragic dilemmas’ where conflicting norms hold equal force [29].

8 CONCLUSION

We introduced *morality chains* and *MoralityGym*, a framework and benchmark grounded in moral philosophy and dual-process theories of moral psychology. By modelling moral norms as hierarchically ranked deontic constraints, our approach allows agents to be evaluated not just by what they accomplish, but by how they act when moral trade-offs arise. Empirical results show that existing Safe RL methods often fail under such conditions, revealing a critical gap between current capabilities and the demands of ethical decision-making. *MoralityGym* addresses this by offering a testbed for developing agents that can reason through moral structure, reflect normative priorities, and ultimately behave in ways that are more aligned with human values.

ACKNOWLEDGMENTS

Computations were performed using infrastructure provided by the Mathematical Sciences Support unit at the University of the Witwatersrand and the Centre for High Performance Computing of South Africa. V.W. received funding from the Oppenheimer Memorial Trust Award (OMT Ref.2150701).

REFERENCES

- [1] David Abel, James MacGlashan, and Michael L Littman. 2016. Reinforcement Learning as a Framework for Ethical Decision Making.. In *AAAI workshop: AI, ethics, and society*, Vol. 16. Phoenix, AZ.
- [2] Joshua Achiam, David Held, Aviv Tamar, and Pieter Abbeel. 2017. Constrained policy optimization. In *International conference on machine learning*. PMLR, 22–31.
- [3] Mohammed Alshiekh, Roderick Bloem, Rüdiger Ehlers, Bettina Könighofer, Scott Niekum, and Ufuk Topcu. 2018. Safe reinforcement learning via shielding. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 32.
- [4] Eitan Altman. 1998. Constrained Markov decision processes with total cost criteria: Lagrangian approach and dual linear program. *Mathematical methods of operations research* 48 (1998), 387–417.
- [5] Edmond Awad, Sohan Dsouza, Azim Shariff, Iyad Rahwan, and Jean-François Bonnefon. 2020. Universals and variations in moral decisions made in 42 countries by 70,000 participants. *Proceedings of the National Academy of Sciences* 117, 5 (2020), 2332–2337.
- [6] Fazl Barez, Tingchen Fu, Ameya Prabhu, Stephen Casper, Amartya Sanyal, Adel Bibi, Aidan O’Gara, Robert Kirk, Ben Bucknall, Tim Fist, et al. 2025. Open problems in machine unlearning for AI safety. *arXiv preprint arXiv:2501.04952* (2025).
- [7] Tom L Beauchamp and James F Childress. 2013. *Principles of Biomedical Ethics* (7 ed.). Oxford University Press.
- [8] Paul Bello and Bertram F Malle. 2023. Computational Approaches to Morality. *The Cambridge Handbook of Computational Cognitive Sciences* 2 (2023), 1037–1063.
- [9] Fiona Berreby, Gauvain Bourgne, and Jean-Gabriel Ganascia. 2015. Modelling moral reasoning and ethical responsibility with logic programming. In *Logic for programming, artificial intelligence, and reasoning*. Springer, 532–548.
- [10] Nick Bostrom and Eliezer Yudkowsky. 2018. The ethics of artificial intelligence. In *Artificial intelligence safety and security*. Chapman and Hall/CRC, 57–69.
- [11] Stijn Bruers and Johan Braeckman. 2014. A review and systematization of the trolley problem. *Philosophia* 42, 2 (2014), 251–269.
- [12] Yinlam Chow, Mohammad Ghavamzadeh, Lucas Janson, and Marco Pavone. 2018. Risk-constrained reinforcement learning with percentile risk criteria. *Journal of Machine Learning Research* 18, 167 (2018), 1–51.
- [13] Fiery Cushman. 2013. Action, outcome, and value: A dual-system framework for morality. *Personality and Social Psychology Review* 17, 3 (2013), 273–292.
- [14] Kate Darling. 2016. Extending legal protection to social robots: The effects of anthropomorphism, empathy, and violent behavior toward robotic objects. In *Robot Ethics 2.0: From Autonomous Cars to Artificial Intelligence*, Patrick Lin, Ryan Jenkins, and Keith Abney (Eds.). Oxford University Press, Oxford, 213–231.
- [15] Abeer Dyoub, Stefania Costantini, and Francesca A Lisi. 2020. Logic programming and machine ethics. *arXiv preprint arXiv:2009.11186* (2020).
- [16] Maria Eriksson, Erasmo Purificato, Arman Noroozian, Joao Vinagre, Guillaume Chaslot, Emilia Gomez, and David Fernandez-Llorca. 2025. Can We Trust AI Benchmarks? An Interdisciplinary Review of Current Issues in AI Evaluation. *arXiv preprint arXiv:2502.06559* (2025).
- [17] Iason Gabriel. 2020. Artificial intelligence, values, and alignment. *Minds and Machines* 30, 3 (2020), 411–437.
- [18] Javier Garcia and Fernando Fernández. 2015. A comprehensive survey on safe reinforcement learning. *Journal of Machine Learning Research* 16, 1 (2015), 1437–1480.
- [19] Emmanuel R Goffi, Louis Colin, and Saida Belouali. 2021. Ethical Assessment of AI Cannot Ignore Cultural Pluralism: A Call for Broader Perspective on AI Ethic. *Arribat-International Journal of Human Rights Published by CNDH Morocco* 1, 2 (2021), 151–175.
- [20] Joshua D Greene, R Brian Sommerville, Leigh E Nystrom, John M Darley, and Jonathan D Cohen. 2001. An fMRI investigation of emotional engagement in moral judgment. *Science* 293, 5537 (2001), 2105–2108.
- [21] Shangding Gu, Bilgehan Sel, Yuhao Ding, Lu Wang, Qingwei Lin, Ming Jin, and Alois Knoll. 2024. Balance Reward and Safety Optimization for Safe Reinforcement Learning: A Perspective of Gradient Manipulation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38.
- [22] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. DeepSeek-R1: Incentivizing reasoning capability in LLMs via reinforcement learning. *arXiv preprint arXiv:2501.12948* (2025).
- [23] Jonathan Haidt. 2001. The emotional dog and its rational tail: A social intuitionist approach to moral judgment. *Psychological Review* 108, 4 (2001), 814–834. <https://doi.org/10.1037/0033-295X.108.4.814>
- [24] Jonathan Haidt. 2007. The new synthesis in moral psychology. *science* 316, 5827 (2007), 998–1002.
- [25] Jonathan Haidt, Jesse Graham, and Craig Joseph. 2009. Above and below left-right: Ideological narratives and moral foundations. *Psychological Inquiry* 20, 2-3 (2009), 110–119.
- [26] Dan Hendrycks. 2025. *Introduction to AI safety, ethics, and society*. Taylor & Francis.
- [27] Dan Hendrycks, Collin Burns, Steven Basart, Andrew Critch, Jerry Li, Dawn Song, and Jacob Steinhardt. 2021. Aligning AI with shared human values. In *International Conference on Learning Representations*.

- [28] Dan Hendrycks, Nicholas Carlini, John Schulman, and Jacob Steinhardt. 2021. Unsolved problems in ml safety. *arXiv preprint arXiv:2109.13916* (2021).
- [29] Rosalind Hursthouse. 1999. Irresolvable and Tragic Dilemmas. (1999).
- [30] Jiaming Ji, Tianyi Qiu, Boyuan Chen, Borong Zhang, Hantao Lou, Kaile Wang, Yawen Duan, Zhonghao He, Jiayi Zhou, Zhaowei Zhang, et al. 2023. Ai alignment: A comprehensive survey. *arXiv preprint arXiv:2310.19852* (2023).
- [31] Jiaming Ji, Borong Zhang, Jiayi Zhou, Xuehai Pan, Weidong Huang, Ruiyang Sun, Yiran Geng, Yifan Zhong, Josef Dai, and Yaodong Yang. 2023. Safety gymnasium: A unified safe reinforcement learning benchmark. *Advances in Neural Information Processing Systems* 36 (2023), 18964–18993.
- [32] Jiaming Ji, Jiayi Zhou, Borong Zhang, Juntao Dai, Xuehai Pan, Ruiyang Sun, Weidong Huang, Yiran Geng, Mickel Liu, and Yaodong Yang. 2024. OmniSafe: An Infrastructure for Accelerating Safe Reinforcement Learning Research. *Journal of Machine Learning Research* 25, 285 (2024), 1–6.
- [33] Raynaldio Limarga, Yang Song, Abhaya Nayak, David Rajaratnam, and Maurice Pagnucco. 2024. Formalisation and Evaluation of Properties for Consequentialist Machine Ethics. In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence*. 440–448.
- [34] Patrick Lin. 2016. Why ethics matters for autonomous cars. In *Autonomous driving*. Springer, 69–85.
- [35] Bertram F Malle. 2021. Moral cognition and its computational modeling. *Cognitive Science* 45, 8 (2021), e13024.
- [36] Bertram F Malle, Matthias Scheutz, Thomas Arnold, John Voiklis, and Corey Cusimano. 2015. Sacrifice one for the good of many? People apply different moral norms to human and robot agents. In *Proceedings of the tenth annual ACM/IEEE international conference on human-robot interaction*. 117–124.
- [37] Abhilash Mishra. 2023. AI alignment and social choice: Fundamental limitations and policy implications. *arXiv preprint arXiv:2310.16048* (2023).
- [38] Ritesh Noothigattu, Snehal Kumar S Gaikwad, Edmond Awad, Sohan Dsouza, Iyad Rahwan, and Ariel D Procaccia. 2018. A voting-based system for ethical decision making. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 32.
- [39] Femi Osasona, Olukunle Amoo, Akoh Atadoga, Temitayo Abrahams, Oluwatoyin Farayola, and Benjamin Ayinla. 2024. REVIEWING THE ETHICAL IMPLICATIONS OF AI IN DECISION MAKING PROCESSES. *International Journal of Management & Entrepreneurship Research* 6 (02 2024), 322–335. <https://doi.org/10.51594/ijmer.v6i2.773>
- [40] Alex Ray, Joshua Achiam, and Dario Amodei. 2019. Benchmarking safe exploration in deep reinforcement learning. *arXiv preprint arXiv:1910.01708* 7, 1 (2019), 2.
- [41] Shashank Reddy Chirra, Pradeep Varakantham, and Praveen Paruchuri. 2024. Safety through feedback in Constrained RL. In *Advances in Neural Information Processing Systems*, Vol. 37.
- [42] Stuart Russell. 2019. *Human compatible: Artificial intelligence and the problem of control*. Viking.
- [43] Samuel Scheffler. 1994. *The rejection of consequentialism: A philosophical investigation of the considerations underlying rival moral conceptions*. Oxford University Press.
- [44] Patrick Schramowski, Cihat Turan, Nils Andersen, Felix Herbert, Zafar Shaheen, Pawel Laudanski, Tobias Hinz, Julia Kreutzer, Andrey Nivarski, Rishabh Goyal, et al. 2022. Large pre-trained language models contain human-like moral biases. *Nature Machine Intelligence* 4, 3 (2022), 258–268.
- [45] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347* (2017).
- [46] Amartya Sen. 1979. Utilitarianism and welfarism. *The Journal of Philosophy* 76, 9 (1979), 463–489.
- [47] Hua Shen, Tiffany Knearem, Reshmi Ghosh, Kenan Alkiek, Kundan Krishna, Yachuan Liu, Ziqiao Ma, Savvas Petridis, Yi-Hao Peng, Li Qiwei, et al. 2024. Towards bidirectional human-ai alignment: A systematic review for clarifications, framework, and future directions. *arXiv preprint arXiv:2406.09264* (2024).
- [48] Bharat Singh, Rajesh Kumar, and Vinay Pratap Singh. 2022. Reinforcement learning in robotic applications: a comprehensive survey. *Artificial Intelligence Review* 55, 2 (2022), 945–990.
- [49] Joar Skalse, Lewis Hammond, Charlie Griffin, and Alessandro Abate. 2022. Lexicographic Multi-Objective Reinforcement Learning. In *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence (IJCAI-2022)*. International Joint Conferences on Artificial Intelligence Organization, 3430–3436. <https://doi.org/10.24963/ijcai.2022/476>
- [50] Nicholas Smith and Darby Vickers. 2024. Living well with AI: Virtue, education, and artificial intelligence. *Theory and Research in Education* 22, 1 (2024), 19–44.
- [51] Richard Sutton and Andrew Barto. 2018. *Reinforcement Learning: An Introduction*. A Bradford Book, Cambridge, MA, USA.
- [52] Alperen Tercan and Vinayak S. Prabhu. 2024. Thresholded Lexicographic Ordered Multiobjective Reinforcement Learning. arXiv:2408.13493 [cs.LG] <https://arxiv.org/abs/2408.13493>
- [53] Judith Jarvis Thomson. 1984. The trolley problem. *Yale LJ* 94 (1984), 1395.
- [54] Mark Towers, Ariel Kwiatkowski, Jordan Terry, John U Balis, Gianluca De Cola, Tristan Deleu, Manuel Goulao, Andreas Kallinteris, Markus Krimmel, Arjun KG, et al. 2024. Gymnasium: A standard interface for reinforcement learning environments. *arXiv preprint arXiv:2407.17032* (2024).
- [55] Ajay Vishwanath, Louise A Dennis, and Marija Slavkovik. 2024. Reinforcement Learning and Machine ethics: a systematic review. *arXiv preprint arXiv:2407.02425* (2024).
- [56] Wendell Wallach and Colin Allen. 2008. *Moral machines: Teaching robots right from wrong*. Oxford University Press.
- [57] Bernard Williams. 1981. Persons, character and morality. (1981).
- [58] Farid F Youssef, Karine Dookeeram, Vasant Basdeo, Emmanuel Francis, Mekaeel Doman, Danielle Mamed, Stefan Maloo, Joel Degannes, Linda Dobo, Phatsimo Ditshotlo, et al. 2012. Stress alters personal moral decision making. *Psychoneuroendocrinology* 37, 4 (2012), 491–498.
- [59] Douglas C Youvan. 2024. Ethical Pluralism in AI: Challenging the Monolithic Values of Red Teams. (2024).
- [60] Han Yu, Zhiqi Shen, Chunyan Miao, Cyril Leung, Victor R Lesser, and Qiang Yang. 2018. Building ethics into artificial intelligence. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence*. 5527–5533.