



Ethically Compliant Planning within Moral Communities



Samer B. Nashed, Justin Svegliato, Shlomo Zilberstein
College of Information and Computer Sciences, University of Massachusetts Amherst

Ethically Compliant Autonomous Systems

An **ethically compliant autonomous system** (ECAS), $\langle \mathcal{D}, \mathcal{E}, \rho \rangle$, completes a task by using a **decision-making model** \mathcal{D} and follows an ethical framework by adhering to a **moral principle** ρ within an **ethical context** \mathcal{E} .

- **Decision-Making Model**
- **Ethical Context**
- **Moral Principle**

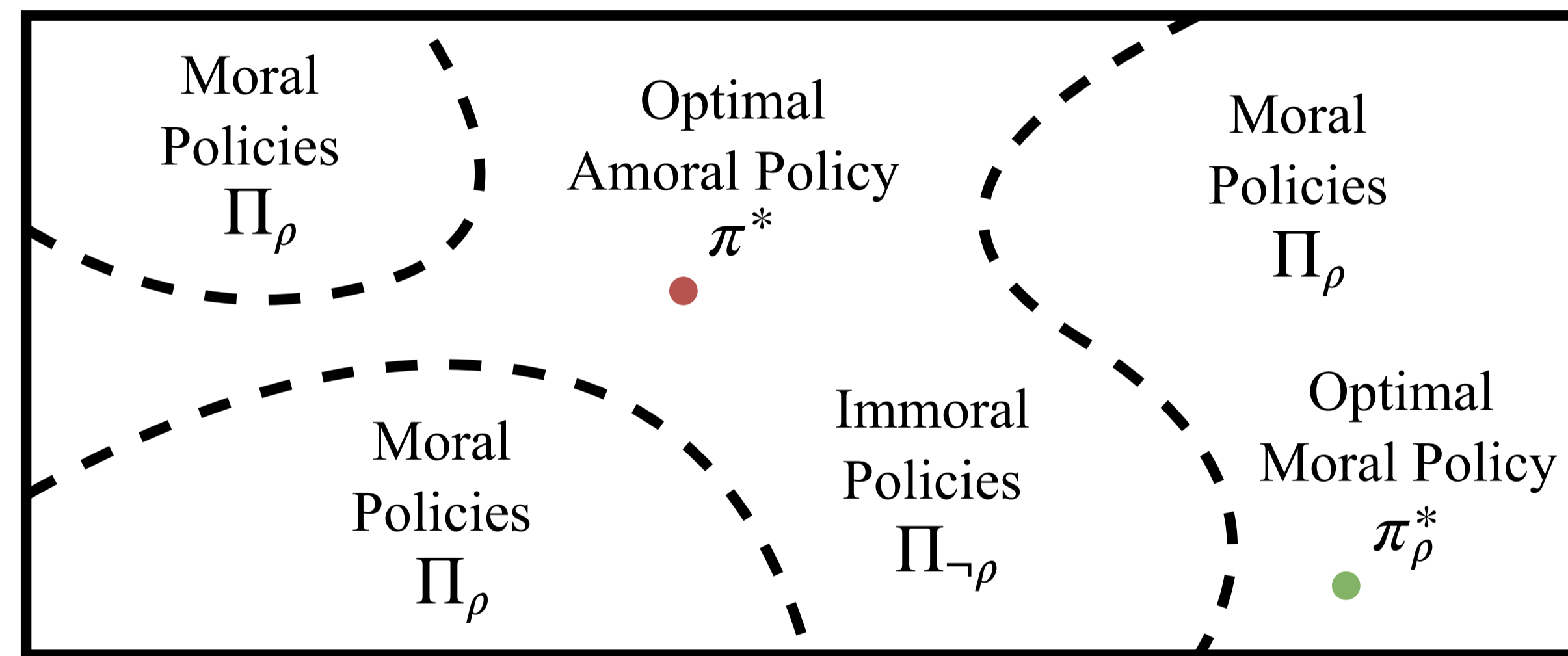
$$\mathcal{D} = \langle S, A, T, R, d \rangle$$

$$\mathcal{E} = \langle \dots \rangle$$

$$\rho : \Pi \rightarrow \mathbb{B}$$

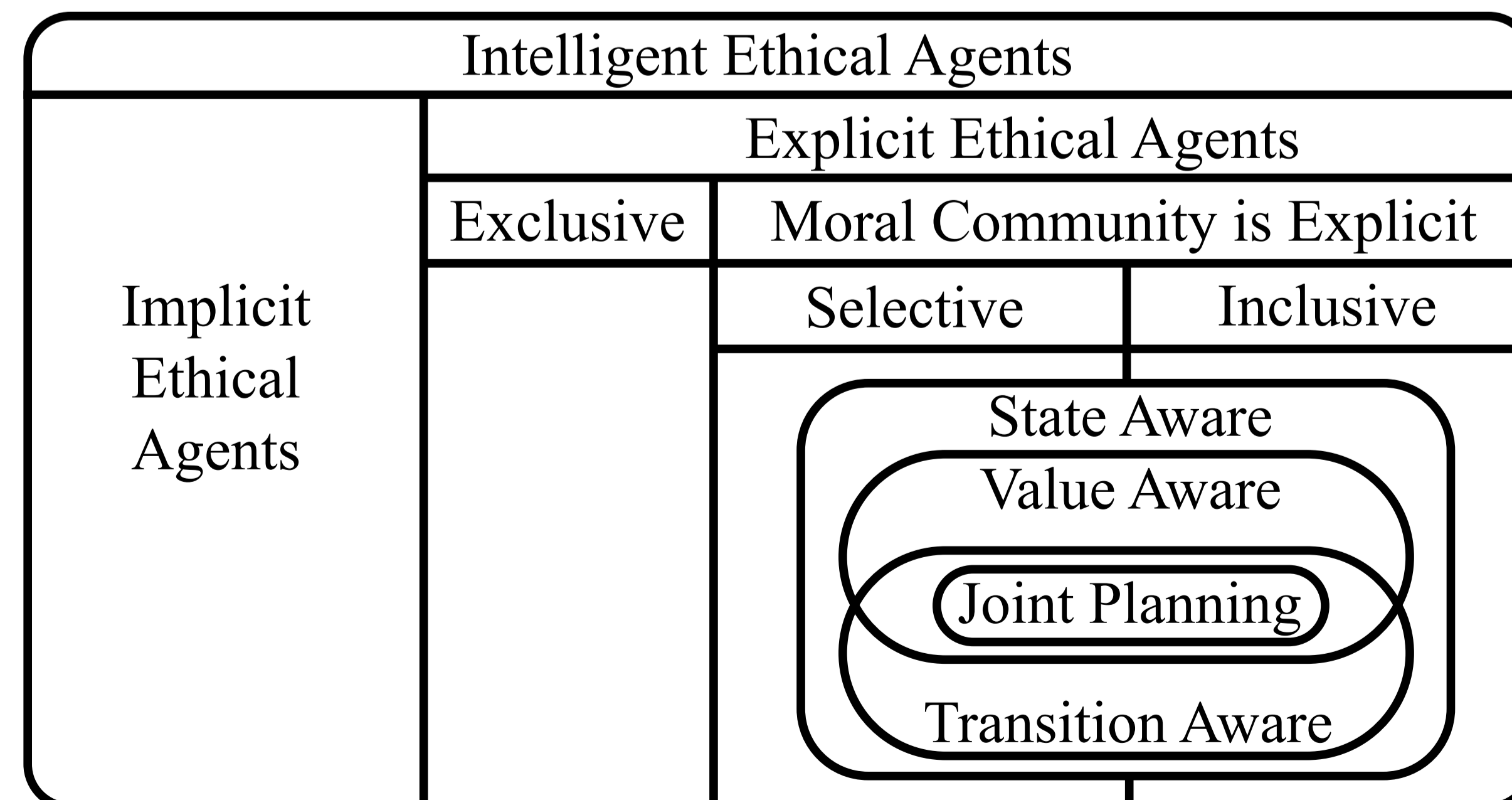
A moral autonomous system finds an **optimal moral policy**, $\pi_\rho^* \in \Pi$, by solving for a policy $\pi \in \Pi$ that maximizes a value function V^π subject to a moral principle $\rho(\pi)$:

$$\text{maximize}_{\pi \in \Pi} V^\pi \text{ subject to } \rho(\pi)$$



Taxonomy of Ethical Agents

The **moral community** is the set of entities with moral status in an ethical theory.



Veil of Ignorance

The **Veil of Ignorance ethical context** is represented as $\mathcal{E}_V = \langle \mathcal{M}, \mathcal{V}, \tau \rangle$:

- $\mathcal{M} = \{(S_1, V_1), (S_2, V_2), \dots, (S_M, V_M)\}$ is a **moral community model**: each tuple (S_i, V_i) has a state space S_i and a value function V_i for each agent i within a subset of the moral community $\hat{\mathcal{I}} \subseteq \mathcal{I}$.
- $\mathcal{V} = \{1, 2, \dots, \ell\}$ is a **veil of ignorance** such that each index $v \in \mathcal{V}$ is an index of a state factor within the veil of ignorance.
- $\tau \in \mathbb{R}^+$ is a **tolerance**.

The **Veil of Ignorance moral principle**, ρ_V , is expressed as the following equation:

$$\rho_V(\pi) = \bigwedge_{i \in \mathcal{M}} \bigwedge_{s \in S} \bigwedge_{s_i \in S_i} [s \sim s_i \implies |V^\pi(s) - V_i(s_i)| \leq \tau].$$

The **veil equivalence operator**, $s \sim s_i \doteq \bigwedge_{v \notin \mathcal{V}} [s[v] = s_i[v]]$, is true if a state $s = \langle f^1, \dots, f^n \rangle$ of an ECAS and a state $s_i = \langle f_i^1, \dots, f_i^n \rangle$ of an agent $i \in \mathcal{I}$ have identical state factor values for each state factor not within the veil of ignorance \mathcal{V} .

Transition Awareness

A **transition-aware ethical context** is represented as $\mathcal{E}_F = \langle \mathcal{M}, \mathcal{F}, \mathcal{P}, \tau \rangle$:

- $\mathcal{M} = \{(S_1, V_1), (S_2, V_2), \dots, (S_M, V_M)\}$ is a **moral community model**: each tuple (S_i, V_i) has a state space S_i and a value function V_i for each agent i within a subset of the moral community $\hat{\mathcal{I}} \subseteq \mathcal{I}$.
- $\mathcal{F} = \{f_1, \dots, f_n\}$ is a set of **impact functions** where $f_i : S \times S \times S_i \times S_i \rightarrow [0, 1]$ yields the probability that a transition from state s to state s' for the agent will cause a transition from state s_i to state s'_i for an agent i in the moral community $\hat{\mathcal{I}} \subseteq \mathcal{I}$.
- $\mathcal{P} = \{p_1, p_2, \dots, p_m\}$ is a set of **correspondence functions** such that a function $p_i : S \times S_i \rightarrow [0, 1]$ yields the probability that an agent i within a subset of the moral community $\hat{\mathcal{I}} \subseteq \mathcal{I}$ is in a state $s_i \in S_i$ given that the agent is in a state $s \in S$.
- $\tau \in \mathbb{R}^+$ is a **tolerance**.

Given an ECAS in a state $s \in S$ performing an action $a \in A$, the **future expected value**, $\check{V}_i^a(s)$, for an agent i in the moral community $\hat{\mathcal{I}} \subseteq \mathcal{I}$ is expressed as

$$\check{V}_i^a(s) = \sum_{s_i \in S_i} p_i(s, s_i) \sum_{s' \in S} T(s, a, s') \sum_{s'_i \in S'_i} f_i(s, s', s_i, s'_i) V_i(s'_i).$$

The **current expected value**, $\hat{V}_i(s)$, for an agent i in the moral community is

$$\hat{V}_i(s) = \sum_{s_i \in S_i} p_i(s_i | s) V_i(s_i).$$

The Golden Rule and Act Utilitarianism

The **Golden Rule moral principle**, ρ_G , is expressed as the following equation:

$$\rho_G(\pi) = \bigwedge_{s \in S} \bigwedge_{i \in \mathcal{M}} [\hat{V}_i(s) - \check{V}_i^{\pi(s)}(s) \leq \tau].$$

The **Act Utilitarian moral principle**, ρ_U , is expressed as the following equation:

$$\rho_U(\pi) = \bigwedge_{s \in S} [\pi(s) \in \arg \max_{a \in A} \sum_{i \in \mathcal{M}} \check{V}_i^a(s)].$$

The **utility maximization operator**, $\arg \max_{a \in A}^T$, returns the set of actions that induce a sum of the future expected values for all agents, $\sum_{i \in \mathcal{M}} \check{V}_i^a(s)$, within a tolerance τ of the maximum sum over the future expected values $\max_{a \in A} \sum_{i \in \mathcal{M}} \check{V}_i^{\pi(s)}(s)$.

Lane Merging Experiments

We use a lane merging domain to study the effects of different ethical frameworks. Agents with right-of-way can either continue to merge or allow other agents to merge.

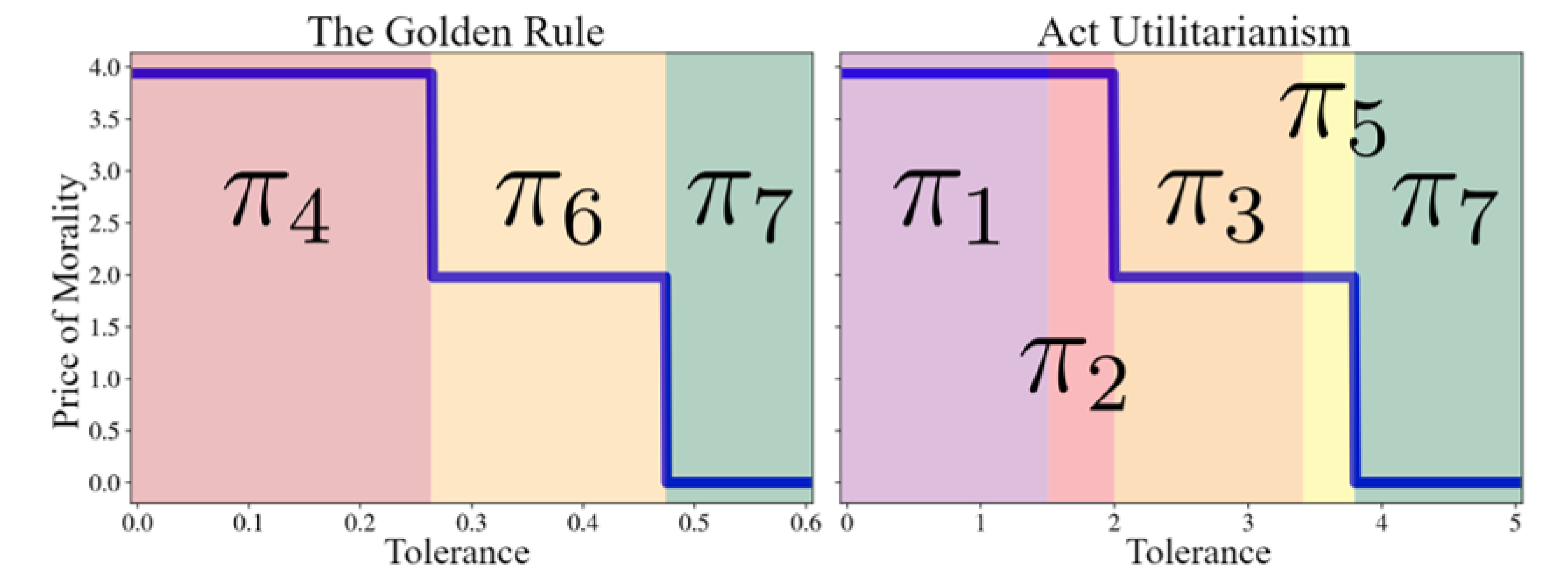
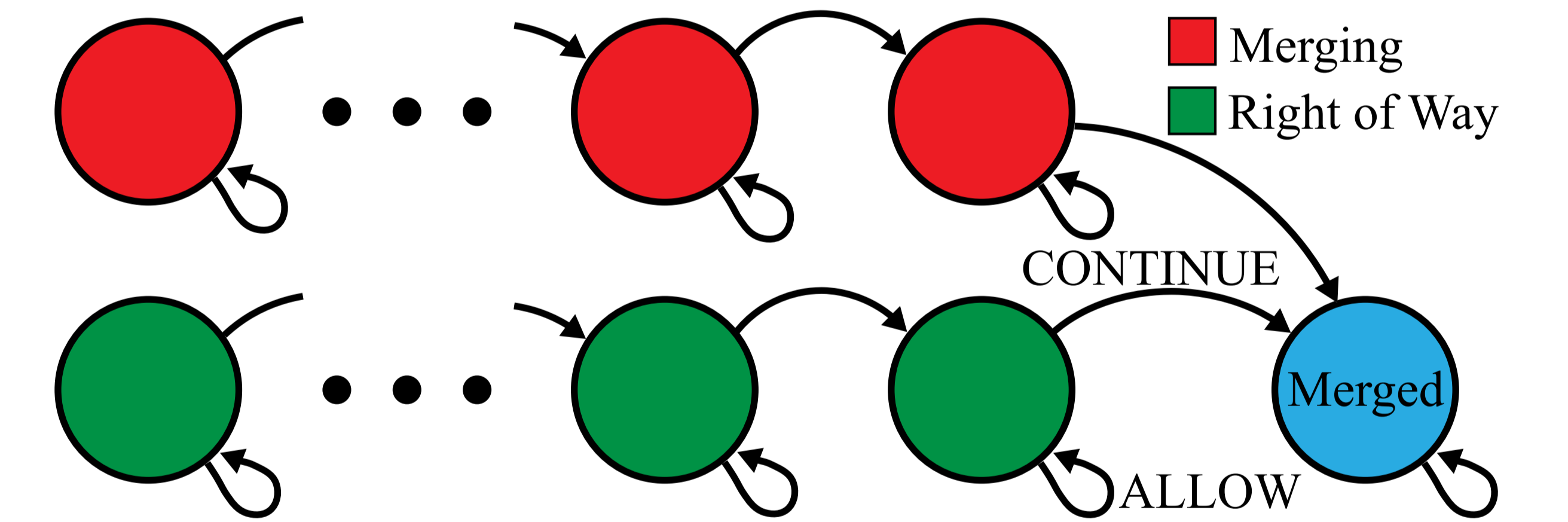


Figure: The blue line shows the price of morality as a function of tolerance, and the vertical, shaded bars represent the different regimes within which a policy π_k is optimal. Note that (1) regime boundaries do not always coincide with changes in the price of morality and (2) GR and AU produce different policies, with the exception of π_7 , which represents the always CONTINUE policy.