

ANALYZING INTENTIONAL BEHAVIOR

IN AUTONOMOUS AGENTS UNDER UNCERTAINTY

Filip Cano Córdoba¹, Samuel Judson², Timos Antonopoulos², Katrine Bjørner³,
Nicholas Shoemaker², Scott J. Shapiro², Ruzica Piskac², Bettina Könighofer¹

¹ Graz University of Technology, ² Yale University, ³ New York University



IJCAI/2023 MACAO

Overview

Accountability: To build trust in autonomous decision-making in uncertain environments it is important to distinguish between *intentional* outcomes, *negligent* desings, and actual *accidents*.

Intention: we propose a definition of intention inspired in Belief-Desire-Intention literature. Taking agents with perfect information as our starting point, we adapt the definition of intention to agents operating in uncertain environments.

Counterfactual reasoning is widely used to study accountability. We ask two types of counterfactual questions:

- *What could the agent have done differently?*
- *How would the agent have behaved in different situations?*

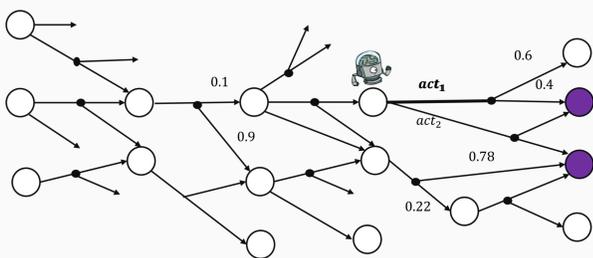
The first question we answer via probabilistic model checking, and the second one we answer by generating counterfactual scenarios.



SCAN ME

Model Setting

- **Environment:** *Markov Decision Processes (MDP)* $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{P})$.
 - **Agent:** Represented as a (deterministic memoryless) policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$.
 - **Intention:** Set of states $\mathcal{S}_I \subset \mathcal{S}$ that the agent intends to reach.
- Probabilistic Model Checking:** What is the probability to reach \mathcal{S}_I ?
- $\mathcal{P}_\pi(\text{Reach}(\mathcal{S}_I), s)$: probability to reach \mathcal{S}_I from $s \in \mathcal{S}$ following policy π .
 - $\mathcal{P}_{\max/\min|\Pi}(\text{Reach}(\mathcal{S}_I), s)$: Max./min. probability, for any policy in $\pi \in \Pi$.



Agency and Intention

Given an agent π at a state $s \in \mathcal{S}$, we define:

- **Scope of Agency** ($\sigma(s)$): Measures the effect of agent's actions on reaching \mathcal{S}_I .

$$\sigma(s) = \mathcal{P}_{\max|\Pi}(\text{Reach}(\mathcal{S}_I), s) - \mathcal{P}_{\min|\Pi}(\text{Reach}(\mathcal{S}_I), s)$$

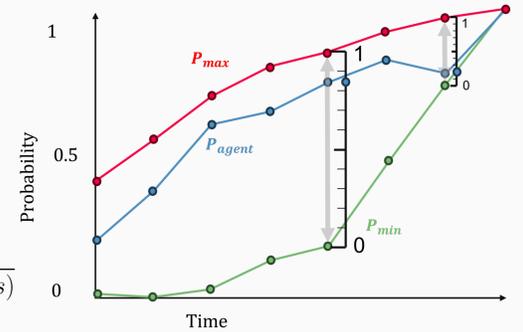
- **Intention-quotient** ($\rho_\pi(s)$): Measures how close π is to being optimal to reach \mathcal{S}_I .

$$\rho_\pi(s) = \frac{\mathcal{P}_\pi(\text{Reach}(\mathcal{S}_I), s) - \mathcal{P}_{\min|\Pi}(\text{Reach}(\mathcal{S}_I), s)}{\mathcal{P}_{\max|\Pi}(\text{Reach}(\mathcal{S}_I), s) - \mathcal{P}_{\min|\Pi}(\text{Reach}(\mathcal{S}_I), s)}$$

For a sequence of states (or trace) $\tau = (s_1, \dots, s_n)$:

Scope of agency ($\bar{\sigma}(\tau)$): Average along a sequence of events (states) τ of the scope of agency.

$$\bar{\sigma}(\tau) = \frac{1}{|\tau|} \sum_{s \in \tau} \sigma(s).$$



Intention-quotient ($\bar{\rho}_\pi(\tau)$): Weighted average along trace τ , weighted by the scope of agency.

$$\bar{\rho}_\pi(\tau) = \frac{1}{\sum_{s \in \tau} \sigma(s)} \sum_{s \in \tau} \sigma(s) \rho_\pi(s)$$

Methodology: Analysis of Intentional Behavior, using Counterfactual Reasoning to Augment Evidence

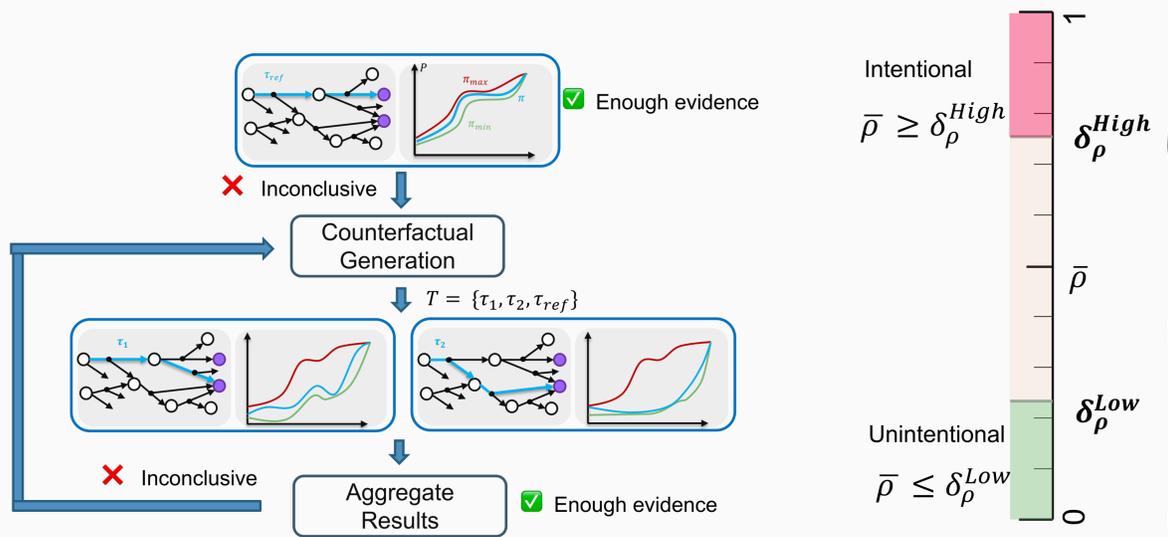
Setting: A factual trace τ that reaches \mathcal{S}_I has happened. We want to analyze whether the agent π reached \mathcal{S}_I intentionally or not. Because of uncertainty, we can only determine *evidence* of intentional behavior. Guiding principle:

- The agent behaves closely to maximizing probability to reach \mathcal{S}_I ,
- The agent could have behaved otherwise.

Thresholds on evidence

- **Agency threshold** (δ_σ). Scope of agency along a trace needs to be larger than the threshold, i.e., $\bar{\sigma}(\tau) \geq \delta_\sigma$. Otherwise, more evidence is required.
- **Intention thresholds** ($\delta_\rho^{\text{High}}, \delta_\rho^{\text{Low}}$).
 - If $\bar{\rho}_\pi(\tau) \geq \delta_\rho^{\text{High}}$, the reaching of \mathcal{S}_I is considered *intentional*.
 - If $\bar{\rho}_\pi(\tau) \leq \delta_\rho^{\text{Low}}$, the reaching of \mathcal{S}_I is considered *unintentional*.
 - Otherwise, more evidence is required.

Counterfactual traces: New counterfactual traces are generated if there is not enough evidence to establish intentionality. Agency and intention-quotient are aggregated along new traces, and compared against thresholds.



Counterfactual Generation

Relevant counterfactuals should be generated with a human in the loop. We propose two semi-automatic generation techniques:

- **Factored MDP.** State space factored into *integral* and *peripheral* state variables $\mathcal{S} = \mathcal{X}_1 \times \dots \times \mathcal{X}_m$. Generate counterfactuals sampling integral variables.
- **Distances on MDPs.** Define a distance notion on states of the MDP, sample traces at close distance.

Discussion & Future Work

Limitations:

- Need for an MDP model of the environment and the agent.
- Probabilistic model checking is costly.
- Agent's beliefs are not taken into account.

Future directions:

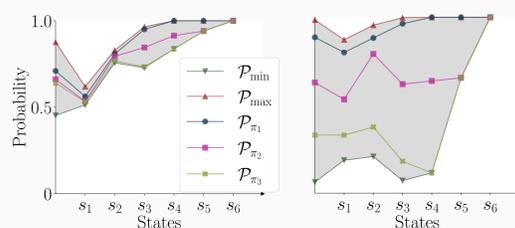
- *General policies.* Extending to policies with memory and non-determinism is feasible, although computationally more expensive.
- *Multi-agent setting.* Considering several agents that interact towards shared or conflicting goals.
- *Time extension.* Longer traces, study intention reconsideration.

Case Study

Example:

- An autonomous car collided with a pedestrian.
- A section of the road was slippery, and there was a truck blocking visibility.

Was the collision intentional?



Comparative analysis of several agents: We built three agents and studied them on the same trace.

- Agent π_1 (●) drives to intentionally hit the pedestrian.
- Agent π_2 (■) drives as fast as possible, caring very little for the safety of the pedestrian.
- Agent π_3 (◆) drives in a safe manner.



Characteristics of the MDP:

- **States:** Position of pedestrian with respect to the car, speed of the car, discretized to integers of m and ms⁻¹.
- **Actions:** Accelerate, brake, coast.
- **Size:** 120k states, 400k transitions.

Generation of counterfactuals: Sample variations of pedestrian behavior, visibility, slippery range, and friction.

