# Causal Explanations for Sequential Decision-Making in Multi-Agent Systems

**Balint Gyevnar,** Cheng Wang, Christopher G. Lucas, Shay B. Cohen, Stefano V. Albrecht

In 23rd International Conference on Autonomous Agents and Multi-Agent Systems



explainable AI · human-centric design · multi-agent systems · cognitive science · natural language processing

**Passenger**                                    **Agent**

Why did you change lanes?

It decreases the time to reach the goal.    } **Teleological**

Why does it decrease the
time to the goal?

Because vehicle 1 was slower than us.

Why was it slower?

It was decelerating and turning right.

What if it hadn't changed
lanes before?                                                **Mechanistic**

We would have gone straight.

# CEMA

**C**ausal **E**xplanations for **M**ulti-**A**gent Systems

## Why multi-agent systems?

Coupled interactions;

Conflicting goals;

Partial observability;

Communication;

## Often difficult to explain, even for humans.

**Critical environments:**

Socially: Others react to our agents and change their behaviour;

Epistemically: Partial observability and shared rules;

Safety: Actions can harm agents/humans/environment;

**Explanations help:**

Explain confusing behaviour;

Highlight occluded information;

Calibrate trust.

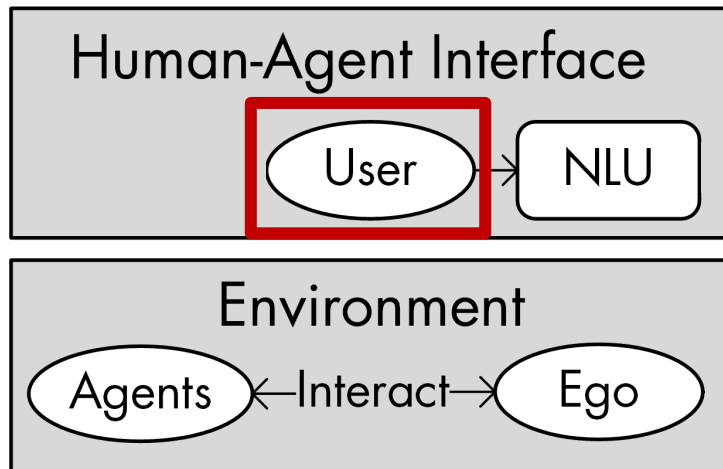**Theory tells us that explanations should be:**
Causal;
Contrastive;
Selected;
Conversational;

[2] Miller, T. (2019). Explanation in artificial intelligence: Insights from the social sciences.
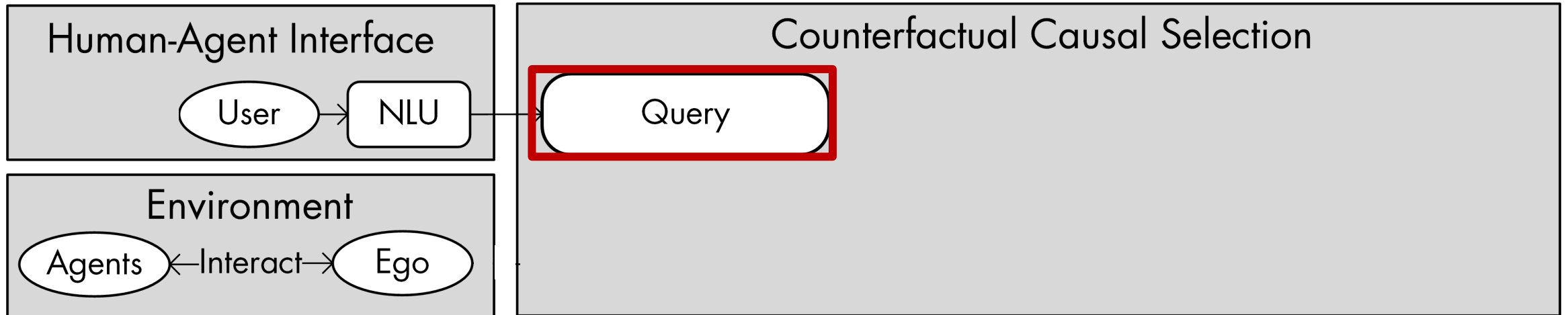
# How does CEMA work?
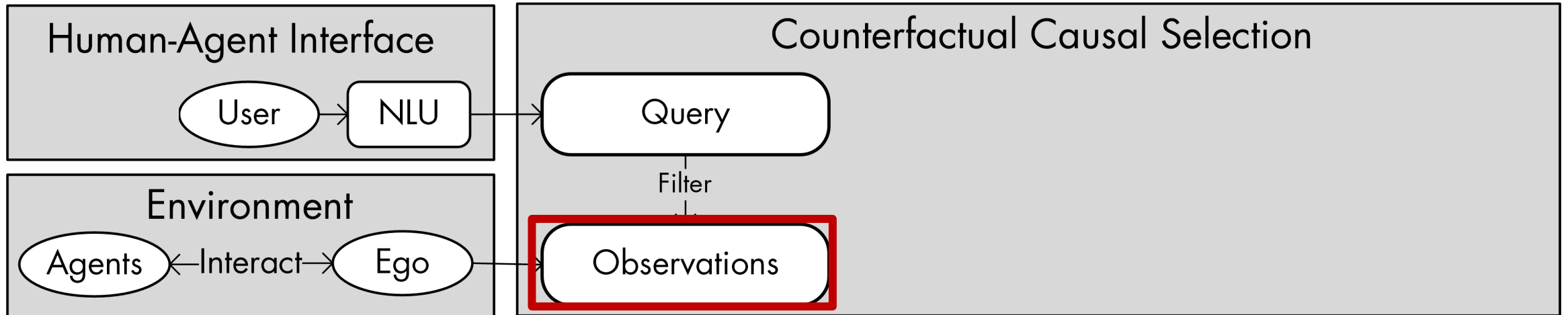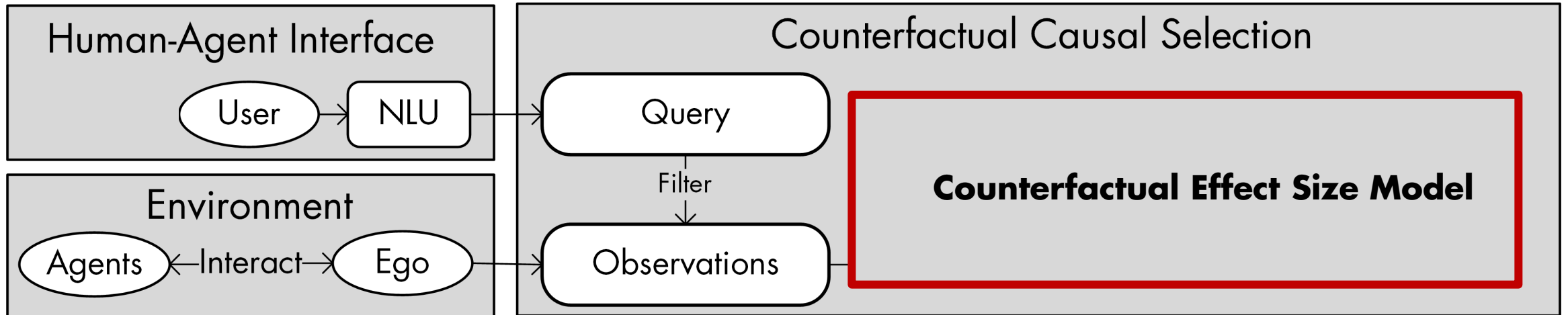
# How does CEMA work?

# How does CEMA work?

# How does CEMA work?

## COUNTERFACTUAL EFFECT SIZE MODEL (CESM)

**Based on how people could select causes to explain;**

**People <u>simulate worlds</u> to select causes:**
Using some prior (cognitive) distribution;
But anchored to observations;

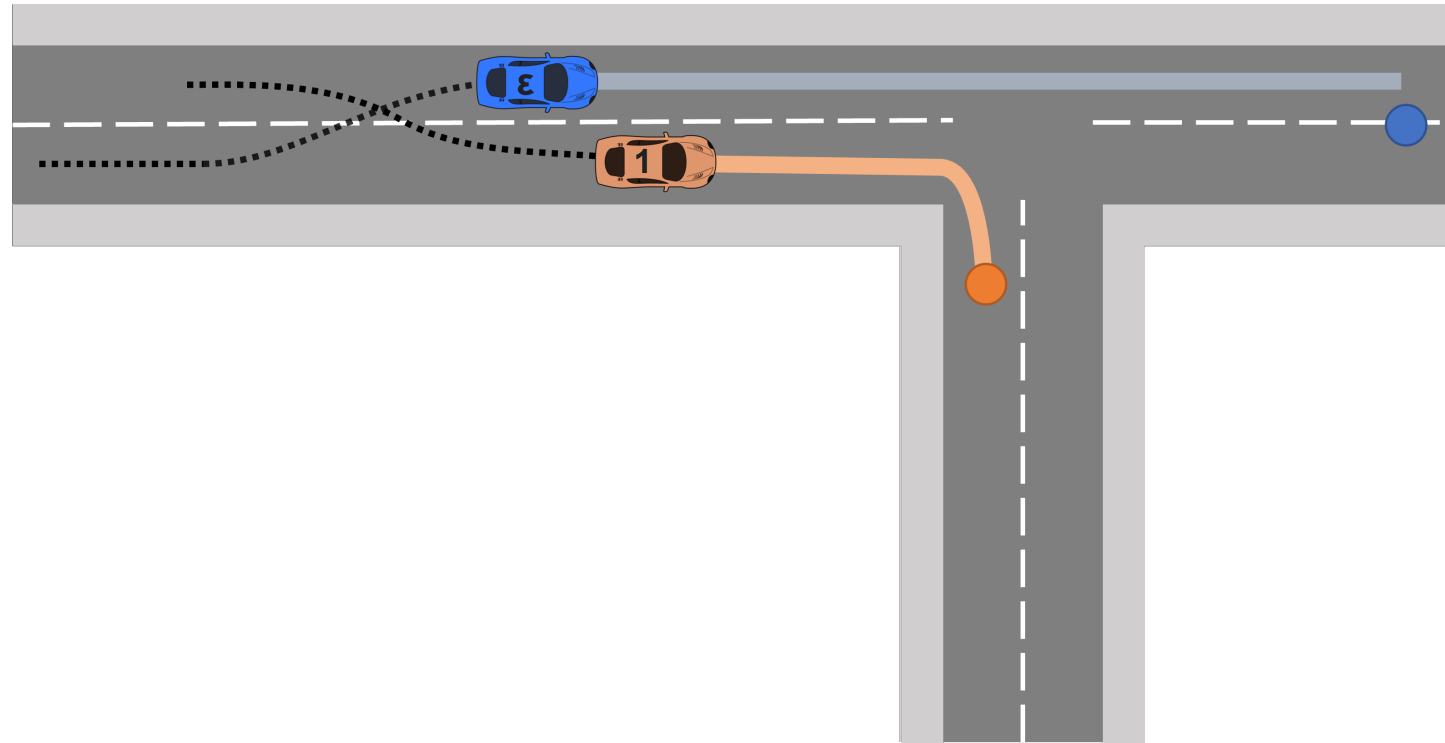**People <u>use correlation</u> to select among causes:**
C caused E if C is highly correlated with E across worlds.

[1] Quillien, T., & Lucas, C. G. (2023, June 8). *Counterfactuals and the Logic of Causal Selection.*
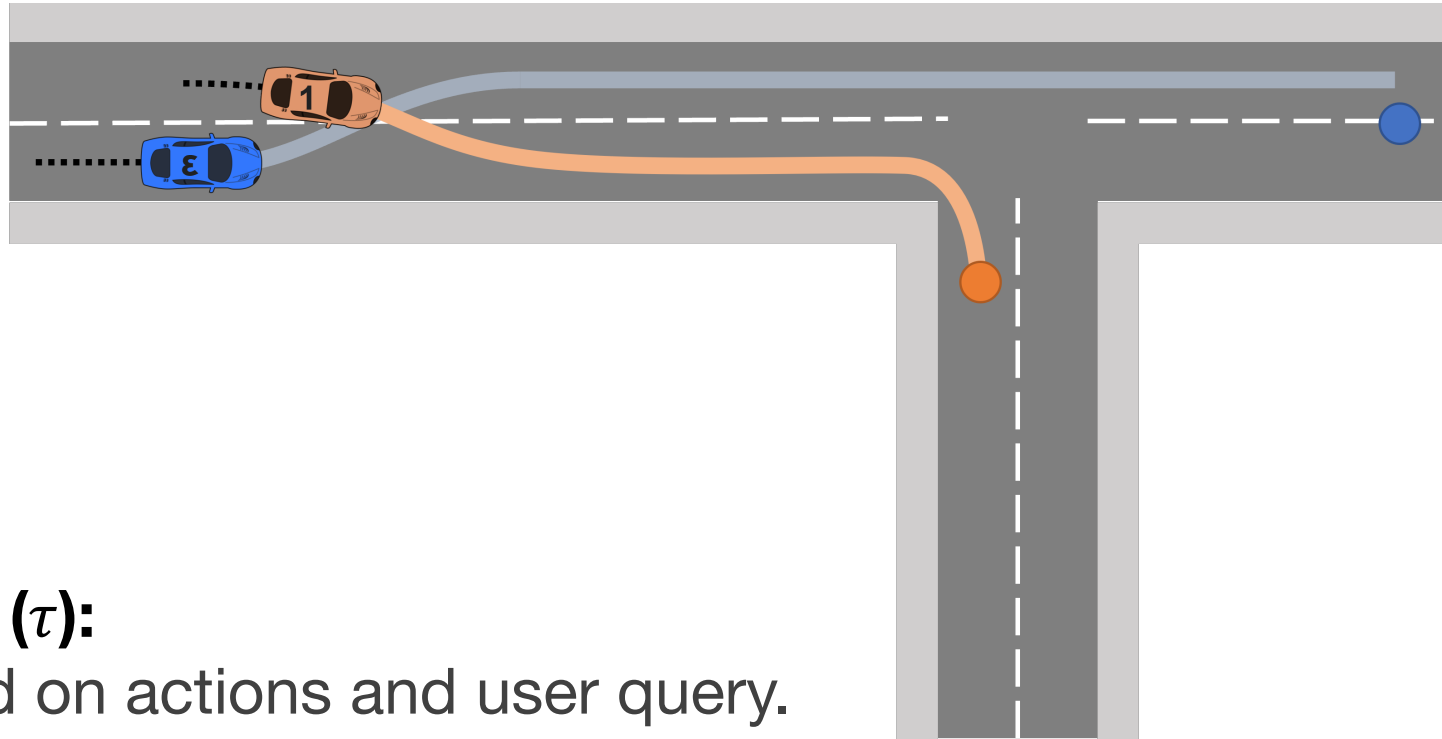
**Counterfactual Effect Size Model in CEMA**

# Rollback → Simulate → Correlate

**Observed trajectory:** $s_{1:t}$



**Rollback** → Simulate → Correlate
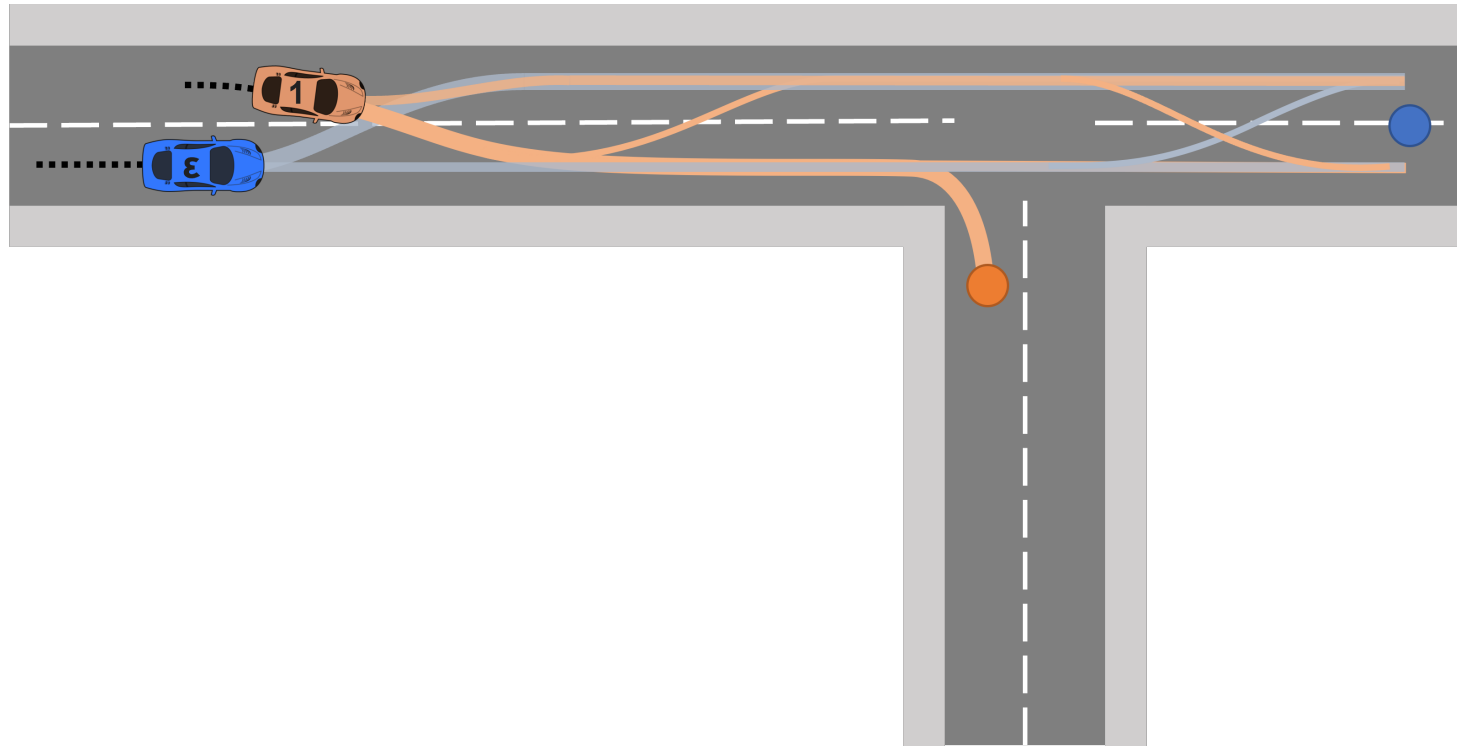
**Rolled back trajectory:** $s_{1:\tau}$



**Rollback time ($\tau$):**
Selected based on actions and user query.

**Rollback** $\rightarrow$ Simulate $\rightarrow$ Correlate

## Simulate (counter)factual worlds



Rollback → **Simulate** → Correlate

**Action presence:**
Lane change (1)

**Rewards:**
Time-to-goal: 7 s
Comfort (jerk): 0.5 m/s$^3$
Collision: No

**Features from trajectory:**
{Decelerate, Turn, Slower, etc…}

Rollback → **Simulate** → Correlate

**Action presence:**
No lane change (0)

**Rewards:**
Time-to-goal: 5 s
Comfort (jerk): 0.1 m/s$^3$
Collision: No

**Features from trajectory:**
{Accelerate, Continue,  Faster, etc…}
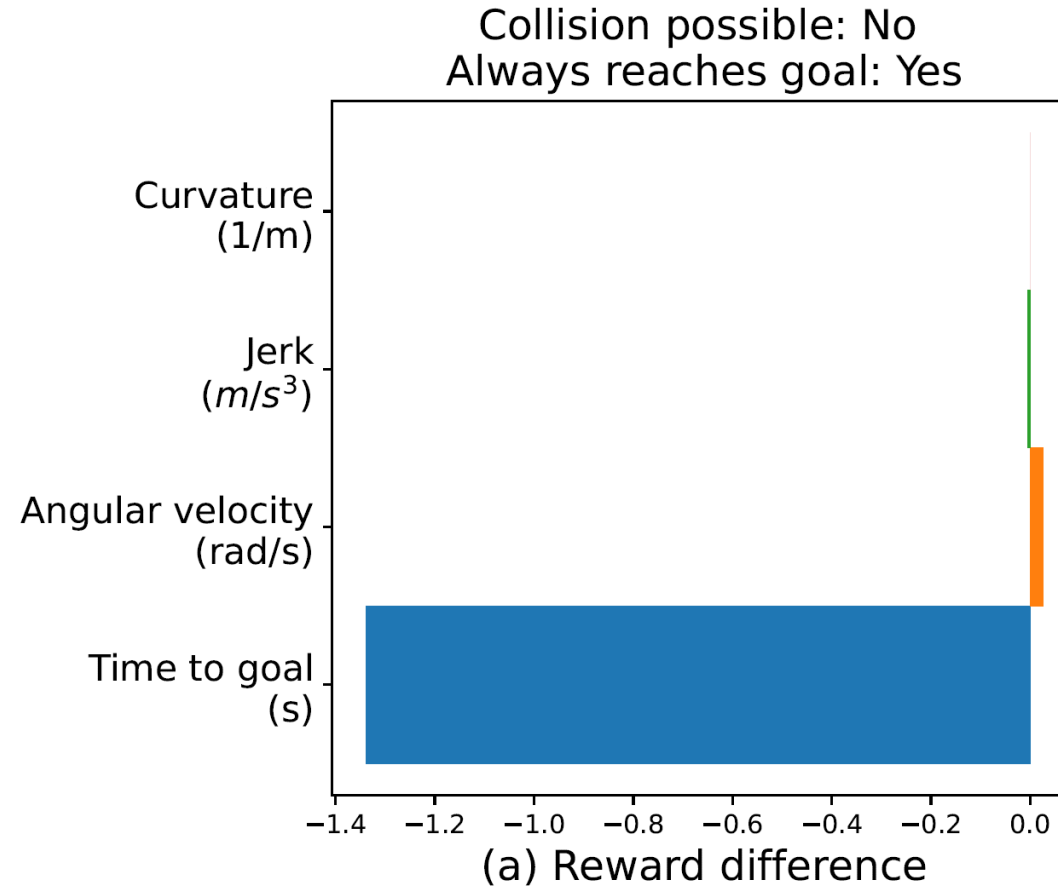
Rollback → **Simulate** → Correlate

**Process for teleological causes;**

**Difference of expected rewards between:**
   Worlds where queried action happened;
   Where queried action did not happen.

Rollback → Sample → **Correlate**

Collision possible: No
Always reaches goal: Yes
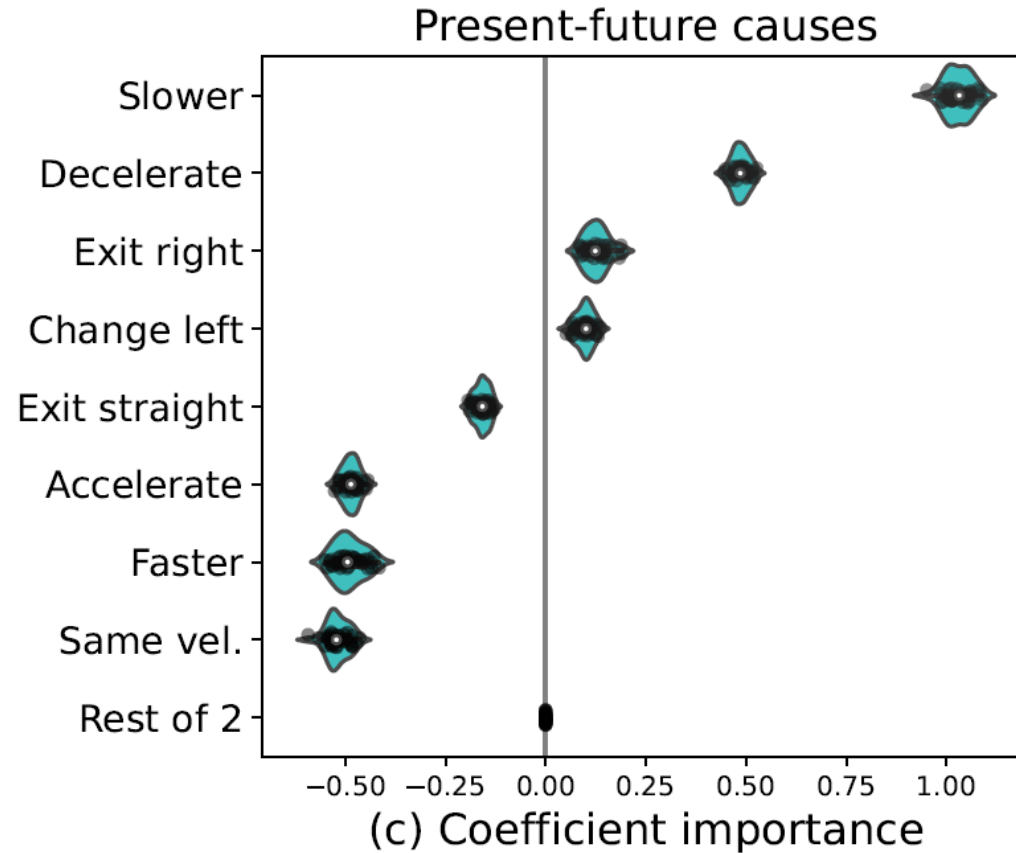
(a) Reward difference

Rollback → Sample → **Correlate**

**Process for mechanistic causes:**
1. Fit interpretable model to trajectory features;
2. Predict presence of queried action;
3. Extract feature importance;

**Counterfactual effect size of features.**
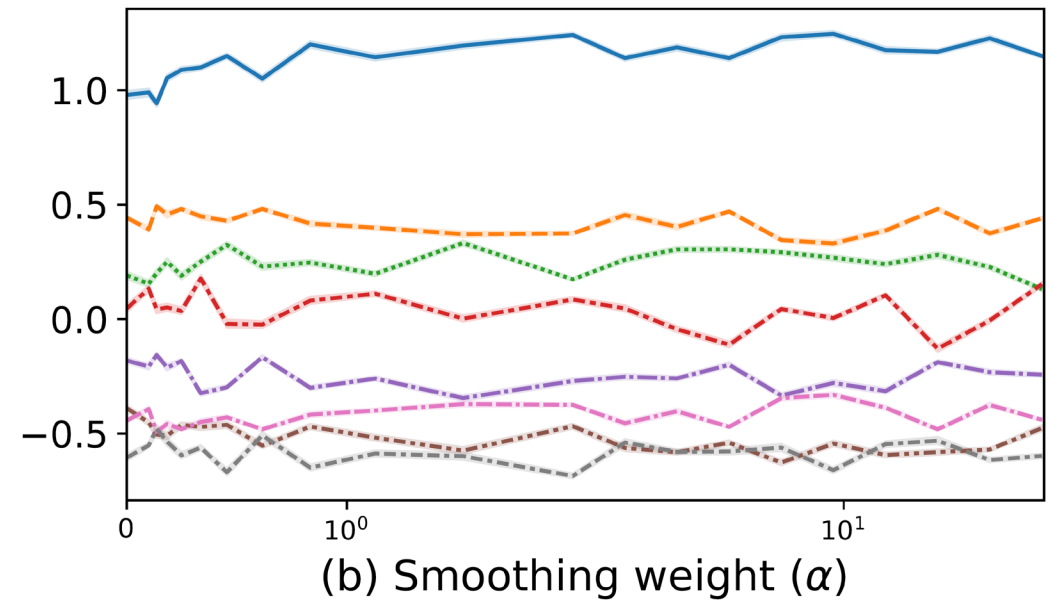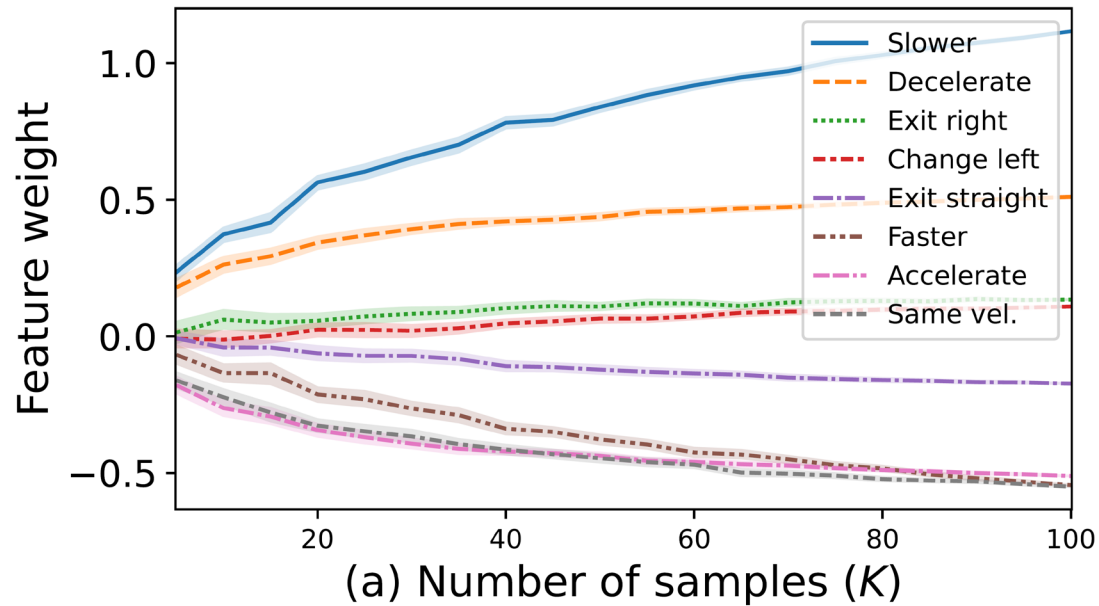
Rollback → Sample → **Correlate**

Present-future causes

(c) Coefficient importance

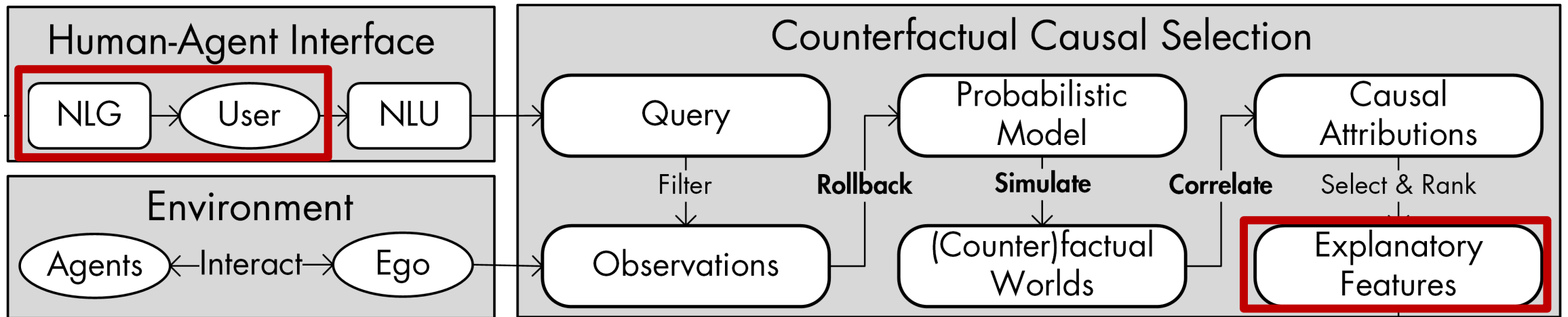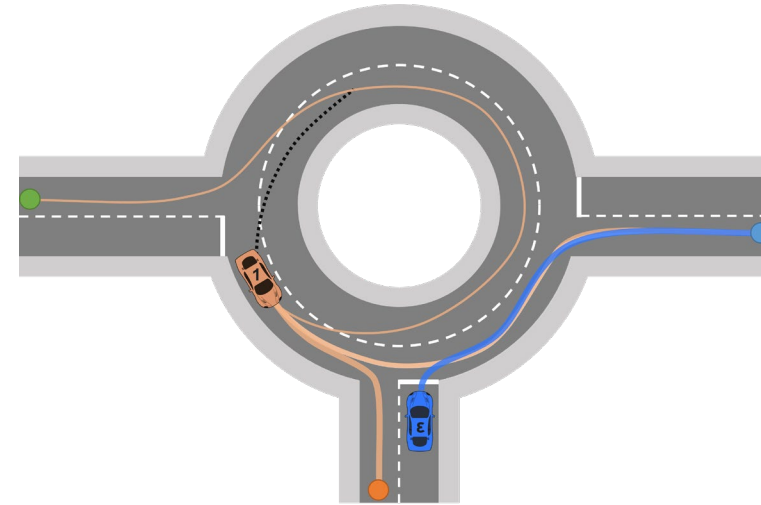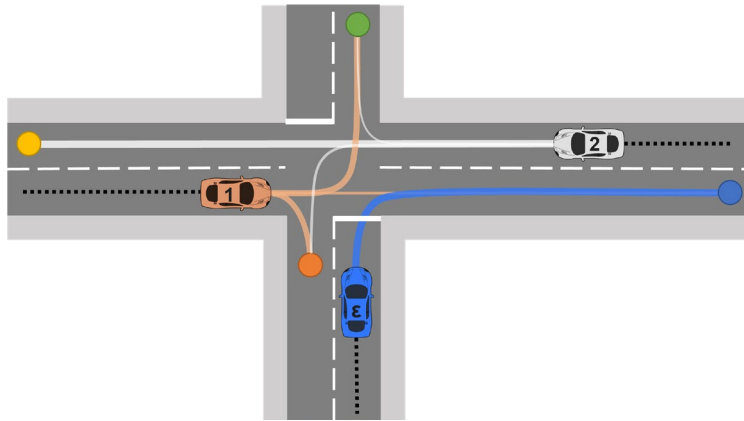Rollback → Sample → **Calculate**

# CEMA is robust



(a) Number of samples ($K$)

(b) Smoothing weight ($\alpha$)

Rollback → Sample → **Calculate**

# FOUR SCENARIOS

## User study:
Generate human-like explanations at least as good as human-written explanations;

## Method:
1. Elicit explanations from people
   (*HEADD: Human Explanations for AD Decisions* dataset [3]);
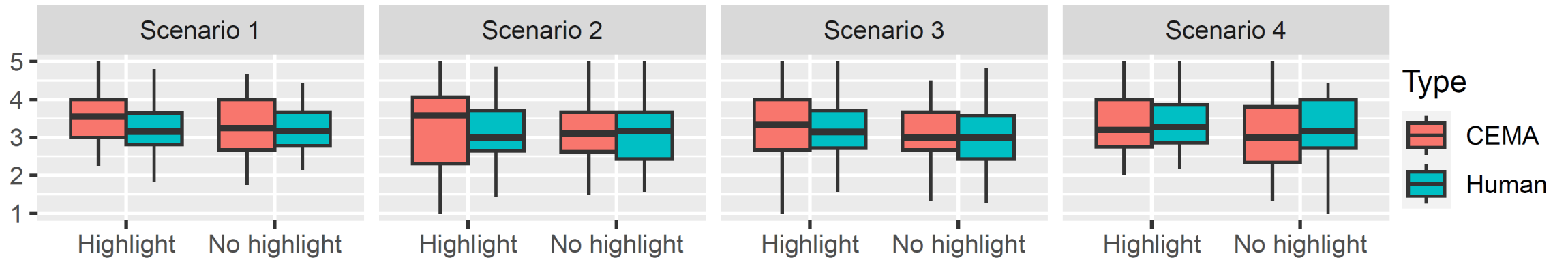2. Compare human explanations to CEMA.

[3] **Gyevnar, B.**, et al. (2024) *People Attribute Purpose to Autonomous Vehicles When Explaining Their Behavior*. (arXiv:2402.10086).

## Independent variables:

Scenarios (1-4)
Explanation type (CEMA/Human)
Highlighting CEMA (Y/N)

➢ **Generally applicable simple (3-step) framework;**

➢ **No explicit assumption on causal structure:**
   No need to model world with DAGs;

➢ **Robust causal selection based on CESM;**

➢ **Works for large number of agents:**
   Tested with up to 20 agents in 4 scenarios.

# Causal Explanations for Sequential Decision-Making in Multi-Agent Systems



https://arxiv.org/abs/2302.10809

## Contributions:

- **CEMA:** General framework for causal explanations of behaviour: <u>Rollback</u> time → <u>Simulate</u> worlds → <u>Correlate</u> variables with outcomes;

- Robustly generated intelligible explanations through natural language.