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

**Balint Gyevnar**
University of Edinburgh
Edinburgh, United Kingdom
balint.gyevnar@ed.ac.uk

**Cheng Wang**
University of Edinburgh
Edinburgh, United Kingdom
cheng.wang@ed.ac.uk

**Christopher G. Lucas**
University of Edinburgh
Edinburgh, United Kingdom
c.lucas@ed.ac.uk

**Shay B. Cohen**
University of Edinburgh
Edinburgh, United Kingdom
scohen@inf.ed.ac.uk

**Stefano V. Albrecht**
University of Edinburgh
Edinburgh, United Kingdom
s.albrecht@ed.ac.uk

## ABSTRACT

We present **CEMA**: **C**ausal **E**xplanations in **M**ulti-**A**gent systems; a framework for creating causal natural language explanations of an agent's decisions in dynamic sequential multi-agent systems to build more trustworthy autonomous agents. Unlike prior work that assumes a fixed causal structure, CEMA only requires a probabilistic model for forward-simulating the state of the system. Using such a model, CEMA simulates counterfactual worlds that identify the salient causes behind the agent's decisions. We evaluate CEMA on the task of motion planning for autonomous driving and test it in diverse simulated scenarios. We show that CEMA correctly and robustly identifies the causes behind the agent's decisions, even when a large number of other agents is present, and show via a user study that CEMA's explanations have a positive effect on participants' trust in autonomous vehicles and are rated as high as high-quality baseline explanations elicited from other participants. We release the collected explanations with annotations as the HEADD dataset.

## KEYWORDS

Explainable AI; human-centric XAI; multi-agent systems; autonomous vehicles; causal explanations; dataset

## 1 INTRODUCTION

Artificial Intelligence (AI) is subject to heightened social and regulatory scrutiny where trust, or a lack thereof, has proven a barrier to public adoption [22], especially in safety-critical systems such as autonomous driving (AD) [18]. This is in part attributed to the inherent lack of transparency of current black box deep learning-based systems [3]. In response, explainable AI (XAI) has gained popularity.
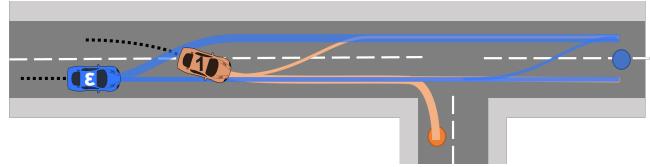
**Figure 1: The autonomous vehicle ($\varepsilon$) is heading to the blue goal. It decided to change lanes after the other vehicle (1) cut in front of it and began to slow down. A passenger asks:** *Why did you change lanes?* **"To decrease the time to reach the goal." [teleological]** *Why was changing lanes faster?* **"Because the other vehicle is slower than us and is decelerating." [mechanistic] – Actual explanations by CEMA with explanation types in brackets. Blue/orange lines illustrate forward simulations using the probabilistic forward model.**

Most XAI methods focus on explanations for supervised learning using tabular or image data [8]. However, these explanations are often purely numeric, and alone have little utility for non-experts who lack domain knowledge to understand the system's internal representations [13]. To address this, XAI is increasingly drawing inspiration from philosophy and the social sciences [28] which has created what we call the subfield of *social XAI*.

An essential part of social XAI is the ability to generate causal explanations. There are several methods for this task [39] and some were proposed for causally explaining sequential decision-making in single-agent systems [9, 40]. However, complex and dynamic multi-agent systems, such as the case with AD, involve tightly coupled interactions among agents where the decisions of any one agent may be difficult to explain even for humans, and there have been few works in XAI addressing this problem. An additional important feature of social XAI is the ability to communicate the extracted causes in the form of intelligible and easy to understand natural language explanations (NLE) as part of a conversational process. A conversation lets users target the pertinent or unclear actions of the agent, while a social XAI system can adjust the user's mental model without excessive cognitive overhead, thereby contributing to more trustworthy interactions with people [11].

To advance the social explainability of multi-agent systems, we introduce a new method called **CEMA**, which stands for **C**ausal **E**xplanations in **M**ulti-**A**gent systems. As illustrated in Figure 1,

CEMA is a social XAI method that generates intelligible causal NLEs about an *ego agent's* decisions in sequential multi-agent systems both in terms of the ego's intrinsic motivations (i.e., teleological explanation) and the actions of other agents in the ego's neighborhood (i.e., mechanistic explanation). At the core of CEMA is a novel causal selection algorithm based on the Counterfactual Effect Size Model [34], which builds on a large body of research into how people select causes for explanations. Instead of creating a specific fixed causal structure, CEMA only relies on a probabilistic model for forward-simulating the joint state of the system, which makes it generally applicable where such models are available. By creating counterfactual simulations of what has occurred, CEMA ranks the salient causes behind the ego's actions based on which causes are most correlated with the ego's actions across counterfactual worlds. Causal selection follows a three-step process:

(1) **Roll back** the current factual state of the system to a previous point in time, such that the actions of the ego that we would like to explain have not yet occurred;

(2) **Simulate** a set of counterfactual worlds from this past time point using a probabilistic forward model of the system;

(3) **Calculate** the counterfactual causal effect size by correlating the ego's actions with changes in its rewards and actions of other agents across counterfactuals.

We evaluate CEMA on AD using diverse simulated driving scenarios from the literature with expert explanations [1], we show that CEMA correctly selects causes of the ego's decisions that are congruent with the expert explanations, even when a large number of agents are present. We show that CEMA is robust to changes in the number of counterfactual simulations and the accuracy of the predictive forward model. We also perform a user study to measure the perceived quality and effects of CEMA's explanations on people. First, we collect a set of high-quality human-written explanations as our baseline. We then show that CEMA's explanations are rated on average at least as high as this baseline while positively affecting participants' trust in AD. In summary, our contributions are:[1]

- CEMA: a framework to generate intelligible causal explanations of the decisions of an ego agent in dynamic multi-agent systems based on the Counterfactual Effect Size Model [34];
- Evaluation of CEMA on motion planning for AD, showing its ability to robustly identify correct causes even when a large number of agents are present;
- HEADD: a dataset of Human Explanations for Autonomous Driving Decisions consisting of human-written explanations with a minimum of 5 unique annotations regarding the causal content and trustworthiness of the explanations [15];
- User study showing CEMA's explanations are ranked at least as high as human explanations and a positive effect of CEMA's explanations on trust in AD.

## 2 BACKGROUND AND RELATED WORK

Causality is a cornerstone of useful human-centric explanations. A common approach for causal selection is to first model the system in the form of a structural causal model (SCM) [32], but this has some drawbacks for complex and dynamically evolving systems.

First, it is challenging to model all causal factors in the system, such as the state, action, or reward influences, while keeping the SCM interpretable and useful for end users. Second, the SCM may grow to intractable sizes depending on the desired coverage of causal factors and the complexity of the system. Third, due to the temporal and non-stationary nature of dynamic systems, an SCM may frequently need to be recomputed to adapt to changes. Thus, existing work has applied SCMs only in simpler single-agent systems where, e.g., the agent is trained with a specific algorithm [27, 29].

In addition, AI models have grown complex enough that generating explanations by "opening the black box", i.e., relying on an understanding of the intrinsic causal properties of the trained model, is often infeasible [42]. Instead, we can rely on the *counterfactual model of causation*, which is a well-understood formulation of causation in philosophical literature [19, 24]. Counterfactual cases uncover causes in relation to the factual case by highlighting events whose absence resulted in the counterfactual case rather than the factual case. Implementing the counterfactual model of causation for complex multi-agent systems is challenging in practice. We rely on Quillien and Lucas [34]'s Counterfactual Effect Size Model which is an empirically validated model to operationalize causal selection based on two assumptions about how humans themselves might select causes for explanations. First, people cognitively simulate counterfactual worlds by sampling from a distribution over possible alternative worlds that are grounded in, i.e., not too different from the factual world. Second, people approximate causal effect sizes by correlating variables (i.e., potential causes) in the world with the presence of an outcome across counterfactual simulations. This means that if we have a probabilistic model for forward-simulating a multi-agent system then we can rank and select the most important causes behind the ego agent's actions by simulating counterfactuals.

Furthermore, how a cause is used for the explanation determines its *explanatory mode*. We consider Aristotle's system as it stood the test of time and is still frequently used in the modern discourse of philosophy of explanations [26]. Aristotle argued for four modes: mechanistic, teleological, material, and formal [16]. The *mechanistic* mode gives an explanation describing the mechanisms of the cause of a change, while the *teleological* mode explains to what end or goal a change has occurred. For example in Figure 1, "other vehicle slowing down" is a mechanistic cause while "reaching goal faster" is a teleological cause behind the decision of the blue autonomous vehicle to change lanes. The material and formal modes stay constant in the systems we study, so we do not consider them.

An increasing body of literature studies the generation of explanations for sequential decision-making. However, most methods focus on deterministic planning in well-defined domains [9]. Prior work in explainable reinforcement learning does address single and multi-agent settings in dynamic systems [33], but causal methods are sparser. Madumal et al. [27] is the first to take a causal approach by building an SCM for the action-influence of agents in model-free RL, while Nashed et al. [29] generates explanation by mapping the algorithmic process of solving a Markov Decision Process into an SCM. Others use surrogate interpretable representations of agents' policies with, e.g., decision trees [38] and programs [41]. We are not aware of methods for social XAI in multi-agent systems.

We use AD for evaluation, where probabilistic models for forward simulating the system are widely available [5]. Goal recognition methods predict other agents' future states [6, 7], while motion planning generates optimal behavior for agents [1, 17]. Social XAI also received some attention in AD. For example, Zhang et al. [44] found that explanations in terms of purely high-level tactical causes (e.g., lane change, turn) had little effect on drivers' trust, therefore, more fine-grained insights are required, e.g., in terms of relative position or acceleration. However, prior methods for social XAI in AD do not consider the sequential nature of decision-making [30], rely on a complex neural model which is impossible to certify for safety [23], or only provide high-level explanations [14].

## 3 CEMA: CAUSAL EXPLANATIONS IN MULTI-AGENT SYSTEMS

We assume that CEMA functions in goal-based sequential multi-agent systems with partial observability, and follow the system definition of Albrecht et al. [1]. Let $\mathcal{I}$ be the set of indexed agents in the environment. At timestep $t \in \mathbb{N}$, each agent $i \in \mathcal{I}$ is in local state $s_t^i \in \mathcal{S}^i$ and receives a local observation $o_t^i \in O^i$ that probabilistically depends on $s_t^i$ through $p(o_t^i \mid s_t^i)$. In addition, agent $i$ selects an action $a_t^i \in \mathcal{A}^i$ in reaction to observations through $p(a_t^i \mid o_{1:t}^i)$, where the notation $o_{a:b}^i$ denotes a tuple for the sequence $(o_a^i, \ldots, o_b^i)$. The joint state of all agents is denoted $s_t \in \mathcal{S}$ where $\mathcal{S} = \times_i \mathcal{S}^i$ and similarly for $o_t \in O$ and $a_t \in \mathcal{A}$. Further, we assume that agent $i$ is aiming to reach a goal $G^i \subset \mathcal{S}^i$ defined as any partial local state description, such as destination coordinates. The goal $G^i$ may not be observable to other agents. If a state sequence $s_{1:t}$ achieves $G^i$ for agent $i$, it receives reward $R^i(s_{1:t}) \in \mathbb{R}^d$ which is a $d$-dimensional vector of reward values where each element in $R^i$ is indexed by a label from a set $\mathcal{R}$ of reward components, such as the time taken to reach the destination. We define the problem of explaining the actions of a particular ego agent $\varepsilon \in \mathcal{I}$ as creating the explanatory function $f : (O^\varepsilon)^* \times (\mathcal{A}^\varepsilon)^* \to \mathcal{H}$ that maps a sequence of local observations and actions to an explanation from a set of possible explanations $\mathcal{H}$. For example, one could define $\mathcal{H} \subset \mathcal{A}^*$, so that an explanation is a partial sequence of actions. We use $\hat{s}_{a:b}$ to indicate that the sequence may contain counterfactual states. We write $s_{x:y} \prec s_{a:b}$ if $s_{x:y}$ is a subsequence of $s_{a:b}$.

We also assume the existence of a probabilistic model that can be used to stochastically forward simulate the system. These are readily available in existing multi-agent literature, for example, in the form of planners or trained reinforcement learning policies [2, 21]. Such probabilistic models define a conditional probability distribution over subsequent joint states of the system given previous observations and actions. We denote this model with $p(\hat{S}_{t+1:n} \mid o_{1:t}^\varepsilon, a_{1:t}^\varepsilon)$, where $n$ is the last timestep. In the case when the local state is fully observable to the ego agent (such as in our evaluation), this model can be replaced with $p(\hat{S}_{t+1:n} \mid s_{1:t}^\varepsilon, a_{1:t}^\varepsilon)$, dropping $a_{1:t}^\varepsilon$ for notational simplicity. Note, that the goals of other agents remain unobservable even under this assumption.

### 3.1 Social XAI Framework

The process of CEMA (Figure 2) begins with the user asking a question about an ego agent $\varepsilon$ and an action they would like explained.
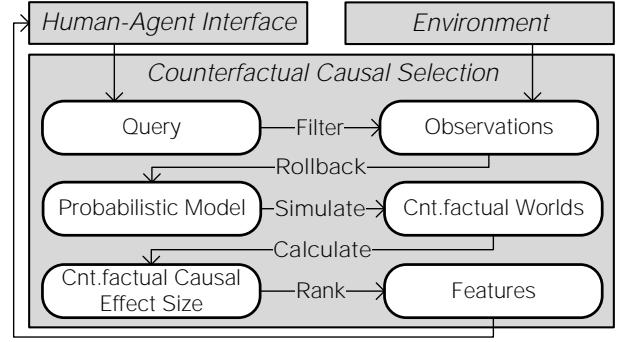


Figure 2: First, irrelevant observations are filtered out based on the query. Second, CEMA rolls back the filtered observations to a previous timestep so that the queried action is erased. From then, CEMA simulates counterfactual worlds to calculate the counterfactual causal effect size for the queried actions, which are used to rank the features of the system.

The question is parsed by an external human-agent interface into a machine-readable *query*, denoted $q$, encoding a description of the state sequence $\hat{s}_{u:v}$ that corresponds to the ego's queried action. Here, $u$ is the start timestep and $v$ is the final timestep of the queried action. Irrelevant states may be then filtered out from the observed states $s_{1:t}$ based on $\hat{s}_{u:v}$. For example, if $\hat{s}_{u:v}$ refers to an action in the past ($v < t$), then we can ignore states after timestep $v$. The queried action $\hat{s}_{u:v}$ need not be a subsequence of $s_{1:t}$, instead it can also be a hypothetical sequence that appears, e.g., in a counterfactual world. This allows the user to ask contrastive questions, for example of the form *"Why did you not do Y instead of X?"* The filtered observations and the query are then passed to the counterfactual causal selection module discussed in detail in Section 3.2.

As the focus of CEMA is to generate intelligible explanations for end users, in this framework explanations are composed from a set of *features* $\mathcal{F}$ which describe semantically meaningful properties of a state and/or action sequence. For discreet $\mathcal{S}$ and $\mathcal{A}$ with inherent interpretations, the set of features might simply equal $\mathcal{S} \cup \mathcal{A}$. For continuous spaces, such as in AD, $\mathcal{F}$ might include a discretized summary of actions, such as average acceleration or distance to the leading vehicle. The set of reward components $\mathcal{R} \subset \mathcal{F}$ are also considered features. For example in autonomous driving, these might be time to destination or presence of collisions. CEMA does not assume anything about the actual meaning or properties of features except that there is some feature function $\phi : \mathcal{S}^* \times \mathcal{A}^* \to \mathcal{F}$ converting a state and action sequence to features. Given the above, for CEMA we define the set of all explanations as $\mathcal{H} = (\mathcal{F} \times \mathbb{R})^*$, so that the output of the counterfactual causal selection process is a subset of features $\mathcal{F}$ with corresponding ranking by counterfactual causal effect size. Finally, the explanation is converted into an NLE and returned to the user via the human-agent interface.

### 3.2 Counterfactual Causal Selection

The counterfactual causal selection process has three main steps. First, it rolls back time before the start timestep $u$ of the queried action, erasing the queried action (Algorithm 1). Second, this rollback

**Algorithm 1** Counterfactual dataset simulation

**Input:** Parsed query $q$; observed joint state sequence $s_{1:t}$.
**Output:** Counterfactual dataset $\mathcal{D} = \{(\hat{s}_{\tau+1:n}^{(k)}, y^{(k)}, r^{(k)})\}_{k=1}^{K}$.

1: $\mathcal{D} \leftarrow \varnothing$.
2: $\tau \leftarrow$ Determine from $s_{1:t}$ assuring that $q.\hat{s}_{u:v}$ is erased.
3: **for** $K$ iterations **do**
4:      Get $\hat{s}_{\tau+1:n} \sim p(\hat{S}_{\tau+1:n} \mid s_{1:\tau})$ via forward simulation.
5:      Determine reward for ego $r \leftarrow R^{\varepsilon}(\hat{s}_{\tau+1:n})$.
6:      Presence of query $y \leftarrow 1$ if $q.\hat{s}_{u:v} < \hat{s}_{\tau+1:n}$ else 0.
7:      $\mathcal{D} \leftarrow \mathcal{D} \cup \{(\hat{s}_{\tau+1:n}, y, r)\}$.
8: **end for**

---

**Algorithm 2** Calculate counterfactual causal effect size

**Input:** Counterfactual dataset $\mathcal{D}$.
**Output:** Mechanistic ($\mathcal{F}^m$) or teleological ($\mathcal{R}^t$) explanation.
*Mechanistic explanation*

1: $\mathcal{F}^m \leftarrow []$.
2: **for** interval end-point $p_j \in P$ **do**
3:      $\mathcal{D}_j \leftarrow$ Slice $\hat{s}_{\tau+1:n}^{(k)} \in \mathcal{D}$ from $p_{j-1}$ to $p_j$ giving $\hat{s}_{p_{j-1}:p_j}^{(k)}$.
4:      $\mathcal{X}, \mathcal{Y} \leftarrow$ convert $\mathcal{D}_j$ to features $\phi(\hat{s}_{p_{j-1}:p_j}^{(k)})$ and targets $y^{(k)}$.
5:      $\mathcal{M} \leftarrow$ Fit an interpretable classifier to $\mathcal{X}$ predicting $\mathcal{Y}$.
6:      $w, I \leftarrow$ Feature importance attributions $w$ of $\mathcal{M}$ indexed in descending order by $I$.
7:      Append $\mathcal{F}_j^m = \{(\mathcal{F}_i, w_i) \mid i \in I\}$ to $\mathcal{F}^m$.
8: **end for**
*Teleological explanation*
9: $\mathcal{X} \leftarrow$ Filter $\mathcal{D}$ by $y^{(k)} = 1$ for match with query.
10: $\mathcal{Y} \leftarrow \mathcal{D} \setminus \mathcal{X}$, all samples not matching the query.
11: $w, I \leftarrow \mathbb{E}_{\mathcal{X}}[r] - \mathbb{E}_{\mathcal{Y}}[r]$ indexed by $I$ in absolute desc. order.
12: $\mathcal{R}^t \leftarrow \{(\mathcal{R}_i, w_i) \mid i \in I\}$.

---

allows CEMA to simulate counterfactual alternatives to the queried action (Algorithm 1). Third, the counterfactual simulations inform us about which features of the system are most important for the queried action to occur and we use this information to calculate the counterfactual causal effect size for both the teleological and the mechanistic explanatory mode presented in Section 2 (Algorithm 2).

Algorithm 1 starts by rolling back the joint state sequence $s_{1:t}$ to a timestep $\tau$, such that $\tau \le u$, resulting in a truncated sequence $s_{1:\tau}$ that assures that the queried action $\hat{s}_{u:v}$ is erased from $s_{1:t}$. The value of $\tau$ can be a fixed distance from $u$ or it can be determined to, for example, correspond to the start of a distinct qualitative change in the ego's behavior prior to $u$. The algorithm then performs $K$ number of forward simulations of the system from time $\tau$ according to the probabilistic model $p(\hat{S}_{\tau+1:n} \mid s_{1:\tau})$. For each simulation, we obtain a sequence of future joint states of the system denoted $\hat{s}_{\tau+1:n}$, determine the reward $r \in \mathbb{R}^d$ for the ego, and whether the queried action $\hat{s}_{u:v}$ of the ego was present in the simulation ($y \in \{0, 1\}$). This process gives a dataset of simulations denoted $\mathcal{D}$.

Algorithm 2 has two parts, one for each mode of explanation.

Mechanistic explanations are formulated in terms of the actions of other agents in the neighborhood of the ego vehicle. Actions of the other agents can have different causal effects on the ego at different times, so we first increase the granularity of explanations by cutting sequences into $|P|$ slices defined by their end-points $P = (p_1, \ldots, p_{|P|})$ with $p_0 = \tau + 1$ assumed implicitly. Each slice is then converted to a set of features using the feature function $\phi$. Following Quillien and Lucas [34], features that co-occur more frequently with the queried action across counterfactuals should be ranked higher as a salient cause by humans. Therefore, for each slice $p_j \in P$ of a counterfactual simulation, Algorithm 2 measures the counterfactual causal effect size of features on the presence of the queried action $y$ by correlating features with the presence of the action across the simulated counterfactuals. For this, an interpretable classifier $\mathcal{M}$ (e.g., logistic regression) is used to predict the presence of the queried action $y$ from the features. The counterfactual causal effect sizes are given by importance attributions for features from $\mathcal{M}$, giving a mechanistic selection and ranking of features $\mathcal{F}_j^m \in \mathcal{H}$.

Teleological explanations are formulated in terms of the intrinsic reward components of the ego agent. For this explanatory mode, the counterfactual simulations inform us how the rewards of the ego, as measured by the reward vector $r \in \mathbb{R}^d$, change depending on the presence $y$ of the queried action of the ego. For binary $y$, this means that Algorithm 2 splits $\mathcal{D}$ into two disjoint sets: one where the queried action was observed ($y = 1$) and one where it was not ($y = 0$). Following the average treatment effect for randomized controlled trials [4] we take the difference between the expected reward vectors of each set, then order the elements of the difference decreasingly by absolute value, giving a teleological ordering of reward components $\mathcal{R}^t \in \mathcal{H}$ by the causal effect of $y$.

## 4 APPLICATION TO MOTION PLANNING

We give a full demonstration of CEMA's capabilities by applying it to the problem of motion planning for AD which is a challenging reasoning task due to the tightly coupled interactions of many agents in a dynamically evolving system [37]. Specifically, we use CEMA to automatically explain the decisions of the Interpretable Goal-based Prediction and Planning (IGP2) system for AD [1]. We give a summary of IGP2 to the extent necessary for the following sections, but for full details please refer to the original paper.

The local state $s^i$ of a vehicle $i$ contains its pose (position and heading), velocity, and acceleration A sequence of temporally adjacent local states is called a trajectory. Local observation $o^i$ contain the local states of nearby traffic participants. Actions $a^i$ set low-level controls such as acceleration and steering, while goals $G^i$ are spatial destinations. Reward components $\mathcal{R}$ are longitudinal and lateral acceleration, presence of collisions, time to reach a destination, and goal completion. IGP2 uses a hierarchy of systems rather than an end-to-end architecture. It defines a set of action sequence templates called *maneuvers* with dynamically generated trajectories for vehicles to follow, including `lane-follow`, `lane-change-{left,right}`, `turn-{left,right}`, `give-way`, and `stop`. Common sequences of maneuvers are then further chained into high-level *macro actions*: `Continue`, `Change-{Left,Right}`, `Exit`, and `Stop`.

IGP2 uses macro actions to predict for each non-ego vehicle $i$ a joint distribution over possible goals and future trajectories given the observed joint local states $s_{1:t}$. Monte Carlo Tree Search (MCTS) is then used to forward simulate the world and obtain driving trajectories for the ego vehicle. In every MCTS simulation, the previously predicted joint goal and trajectory distribution is used to

**Table 1: Binary features $\mathcal{F}$ to describe the fundamental motions and high-level actions of vehicles (including ego). For continuous values, the mean value is calculated along the length of the trajectory and thresholded with small value $\delta$.**

| Feature | Calculation | Explanation |
|---|---|---|
| *Acceleration* | $a^i > \delta_a$ | Accelerate |
| | $a^i < -\delta_a$ | Decelerate |
| | $a^i \in [-\delta_a, \delta_a]$ | Maintain velocity |
| *Relative* | $v^i - v^\varepsilon > \delta_v$ | Faster than ego |
| *speed* | $v^i - v^\varepsilon < -\delta_v$ | Slower than ego |
| | $v^i - v^\varepsilon \in [-\delta_v, \delta_v]$ | Same speed as ego |
| *Stop* | $v^i \in [0, \delta_s]$ | Does it stop |
| *Maneuver* | One-hot encode | Longest maneuver |
| *Macro Action* | One-hot encode | Longest macro action |

randomly sample a goal and corresponding trajectory for each non-ego vehicle. MCTS generates a trajectory for the ego in a simulation by sequentially choosing macro actions based on backpropagated preference values (i.e., $Q$-values) until the ego reaches its goal.

### 4.1 Implementing CEMA

We define our set of features $\mathcal{F}$ in Table 1, which were chosen to describe both fundamental motions and high-level maneuvers of all vehicles including the ego. Features average along the length of the trajectory and may encounter the issue that at one timestep they have a positive causal effect, while at a later timestep, they have a negative causal effect, resulting in aggregate zero causal effect. The slicing operation in Algorithm 2 assures that this issue is avoided.

To focus on causal selection and avoid the ambiguities of natural language, we hand-code each query $q$ to contain a description of the queried subsequence $\hat{s}_{u:v}$ given as a subset of features from $\mathcal{F}$. For natural language generation, we use a deterministic realization engine called SimpleNLG [12], which generates a grammatically correct English sentence from a content specification, e.g., subject and verb. This a better fit than neural generation algorithms, due to a lack of annotated data and hallucinations in neural models.

Since IGP2 can assign to some (reachable) goals and trajectories near-zero probabilities, we use additive smoothing – detailed in Appendix A.3 – with parameter $\alpha$ to make sure every goal and trajectory can be sampled for the non-ego vehicles. We then generate two datasets with Algorithm 1. For teleological explanations, we set $\tau = u$, rolling back time just before the queried action of the ego. This is because teleological explanations are determined by the MCTS reward components which only depend on the ego's present and future actions. For mechanistic explanations, we set $\tau$ to the start time of the last action prior to $u$, erasing both the queried action of the ego and the action that came before it. For slicing the trajectories in Algorithm 2, we set $P \leftarrow (u, n)$ which slices the trajectory $\hat{s}_{\tau+1:n}$ into a past $\hat{s}_{\tau+1:u}$ and present-future $\hat{s}_{u:n}$ subsequence in reference to the start of the ego's queried action.

We use feature weights from logistic regression with K-fold cross-validation to determine feature importance values. We found logistic regression to work best as it is simple, inherently interpretable, and all features are binary so their scale does not affect the importance values.

## 5 COMPUTATIONAL EVALUATION

We evaluate CEMA on the four scenarios (S1–S4) used by Albrecht et al. [1]. The scenarios are shown in Figure 3 with expert explanations of the ego's behavior by Albrecht et al. [1] In line with our focus on social XAI, we test CEMA on many user queries regarding different ego agents and behaviors, and the generated outputs of CEMA are presented through five simulated conversations (Table 2), highlighting CEMA's ability to correctly identify the causes behind each queried action. For all queries, we simulate $K = 100$ counterfactual worlds with a smoothing weight $\alpha = 0.1$. Further details about the experimental setup are given in Appendix B. We focus on S1 for presentation, but all results are confirmed across all scenarios and all presented in Appendix C. We show that:
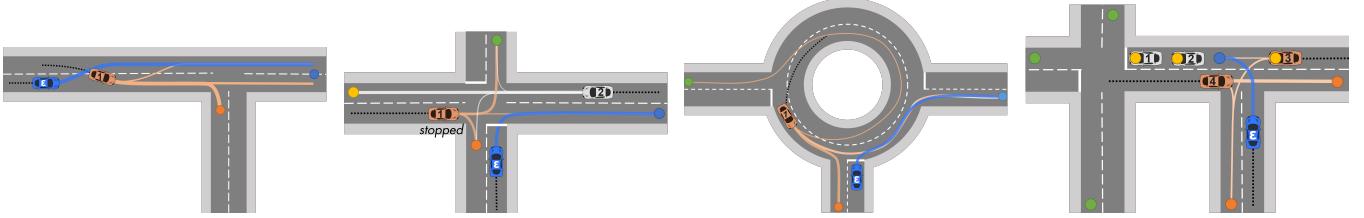
(1) CEMA correctly finds and ranks the relevant causes of the ego's actions that are congruent with expert explanations;
(2) It correctly identifies the relevant causes even in the presence of a large number of agents;
(3) The causal selection process is robust to changes in the sampling size $K$ and the accuracy of the probabilistic model.

### 5.1 Correctness of Causal Selection

As shown in Table 2, CEMA correctly selects causes which are congruent with the expert explanations of Albrecht et al. [1].

In conversation S1-A, the causes behind the factual lane change of the ego are queried. The top plot in Figure 4 shows that CEMA correctly finds that a decrease in time-to-goal is the most significant teleological cause. As the bottom plot in Figure 4 shows, CEMA correctly identifies that the non-ego slowing down is a mechanistic cause of the ego's lane change. CEMA also determines that this slowing down is due to the non-ego vehicle decelerating in order to turn right. The middle plot of Figure 4 confirms that the initially faster non-ego vehicle cutting in front of the ego is also a mechanistic cause of the ego's lane change. This shows the importance of slicing the trajectories into segments as CEMA produces more fine-grained causes that focus on action in a particular time interval.

In conversations S1-B to S4, we also see that CEMA correctly identifies causes for contrastive questions – for example, *"Why aren't you going straight?"* – in which the user asks about an alternative action (i.e., foil) that the ego could have done as opposed to the factual observed actions (i.e., fact). Leveraging the counterfactual simulations, CEMA contrasts the simulations containing the foil to simulations containing the fact and derives the appropriate teleological causes. CEMA delivers consistent explanations even when queries target the same action but are phrased differently. For example, *"Why will you change lanes?"* is a direct question, while *"Why aren't you going straight?"* is contrastive, yet they both refer to the same changing lane action of the ego and CEMA finds consistent causes for both queries. In S4, CEMA correctly finds that the stopping of non-ego 3 is the most relevant cause behind the ego's early merging behavior and it also finds other intuitive causes. For example, the vehicle at the front of the waiting line of cars is

**(S1)** The **non-ego** in front of the **ego** changes lanes and begins to slow down. This is indicative of its intention to turn right at the junction. To avoid being slowed down, the ego decides to change lanes as it is heading straight.

**(S2)** The **ego** is turning right but must give way. It observes the **vehicle** on the left stopping. This is only rational if it is trying to turn left and is giving way for the oncoming **vehicle**. The **ego** can use this to enter the road earlier.

**(S3)** The **ego** observes the **non-ego** changing lanes to the right. This is only rational if the **non-ego** is leaving the roundabout at the next exit. The **ego** can therefore enter the roundabout faster without waiting to give way.

**(S4)** **Non-ego 3** is slowing down to stop. Once **non-ego 4** drives past as indicated by its maintained high speed, the stopping of **non-ego 3** stays rational only if it is to allow the **ego** to merge without waiting for **non-ego 4** to pass.

**Figure 3: The four scenarios used for evaluation based on Albrecht et al. [1]. Colored circles are goals. Solid lines are predicted trajectories of non-egos with thickness corresponding to predicted probability. Black dotted lines are observations.**

stopped. Would this vehicle move, the waiting line of cars would begin moving and non-ego 3 could not allow the ego to merge.

CEMA can also correctly find the relevant causes even when a large number of agents are present. For this, we greatly increase the number of agents in all scenarios and rerun CEMA. For example, we extend S1, adding two extra lanes to the east-west road and increasing the number of agents to 20. This gave 180 features, most of which had no causal influence on the ego, but CEMA could still identify the most important causes as in the original scenario.

## 5.2 Robustness of Causal Selection

We demonstrate robustness to changes in **(a)** the sampling size $K$, and **(b)** in the accuracy of the probabilistic simulation model, to show that correct explanations are generated even when sampling is limited by resources and that our system works with prediction algorithms of varying performance. For size robustness, we randomly sample a dataset of $K \in \{5, 10, \ldots, 100\}$ sequences 50 times and calculate the causal attributions for each dataset. For robustness, we interpolate between the true predicted and uniformly distributed behaviors by increasing the smoothing strength $\alpha$ on a log scale.

The top plot in Figure 5 shows the evolution of causal attributions as we increase $K$ in S1. We see that CEMA becomes increasingly confident in its attributions as $K$ increases, while confidence intervals remain tight. Even with few samples, CEMA identifies causes correctly. The bottom plot of Figure 5 shows how causal attributions change as $\alpha$ increases which corresponds to increasing uncertainty in behavior predictions. We see that feature importance values are little affected by changes in the sample distributions as they fluctuate around the same values. Similar patterns are observed across scenarios, which demonstrates that CEMA is robust to changes in both the sampling size and the accuracy of external predictions.

## 6 USER STUDY

So far, we have focused on the technical details of CEMA. Ultimately, however, the primary target of CEMA is non-expert end users, so we must evaluate the quality of CEMA's explanations and their various effects on humans with actual participants via a user study. We aim to answer the following research questions:

(1) How do people perceive the quality of CEMA's explanations as compared to a human baseline?
(2) What are the effects of CEMA's explanations on people's trust in autonomous vehicles?

We used Prolific to recruit participants from the USA whose first language is English. As most people have not had first-hand experience with autonomous vehicles (AV), we engaged them via animated videos of the scenarios. We design two surveys and summarize our methodology below with full details in Appendix D.

In the first survey (N=54; Male=25, Female=29), participants were asked to describe and explain in their own words the behavior of the AV in all four scenarios. This gave 408 explanations across scenarios, of which we excluded 26 vacuous responses (e.g., *"I don't know"*, *"None", etc.*), and annotated the remaining explanations with a different set of participants regarding their causal content, overall quality and complexity, and trustworthiness. We release an extended version of this annotated dataset of natural language explanations, called the *Human Explanations for Autonomous Driving Decisions (HEADD)* dataset, containing 14 scenarios with several agents and environmental elements, including occlusions, pedestrians, and 1308 explanations. We collected explanations as we are not aware of any reproducible and publicly available methods for AD that would allow for a meaningful comparison to CEMA's explanations.[2] Comparing against a human baseline is also a better fit for CEMA as its explanations are intended to have low cognitive overhead and be easy to understand. In contrast, more complex expert explanations would likely be less effective for end users [11].

In the second survey (N=200; M=99, F=101), we designed two tasks, one for each research question. First, to measure the quality of explanations, we asked participants to rate a random sample of 10 explanations from a set of 30 explanations (5 from CEMA and 25 from HEADD with the highest quality ratings) for each scenario on a 5-point Likert scale. With 50% chance, we highlighted

---

[2]We explored ChatGPT as a baseline, but it was inadequate as its responses were very inconsistent and only sometimes correct (see Appendix D.4 for details).

**Table 2: Actual responses of CEMA to queries (*in italic*). (S1-A)** The passenger sees on an onboard display, that the ego is planning to change to the left lane. They find this unexpected and inquire. **(S1-B)** The passenger observes the ego changing lanes and asks for the reasons behind the maneuver while it is ongoing. **(S2)** The passenger observed that the ego had entered the junction without stopping to give way, which could be dangerous. **(S3)** The passenger sees that the ego will turn right at the roundabout without stopping despite the oncoming vehicle. **(S4)** Once non-ego 4 has passed, the ego merges onto the main road which makes the passenger wonder why the ego did not give way to non-ego 3.

| Scenario 1 (S1-A) | Scenario 1 (S1-B) |
|---|---|
| *Why will you change lanes?* It will decrease the time to goal. *Why does it decrease the time to the goal?* Because vehicle 1 will be slower than us. *Why will it be slower?* It will decelerate and turn right. *What if it hadn't changed lanes before?* We would've gone straight. | *Why aren't you going straight?* It would increase the time to the goal. *Why is it slower to go straight?* Because vehicle 1 is slowing down. *Any other reasons?* It is slower than us. *What will you do after this?* We will go straight and accelerate. |
| **Scenario 2 (S2)** | **Scenario 3 (S3)** |
| *Why did you not stop to give way?* It would've increased the time to the goal. *Couldn't it have caused a collision?* A collision would not have occurred. *How did you know?* Vehicle 1 was turning left and stopped. *What if it went straight?* We would've given way and slowed down. | *What will you do at the roundabout?* We will turn right and accelerate. *Would not stopping lead to a collision?* Not stopping doesn't cause a collision, but stopping increases the time to reach the goal. *How do you know we won't collide with the oncoming car?* It has been changing lanes right and is turning right. |
| **Scenario 4 (S4)** | |
| *Why are you not stopping to give way?* Stopping and giving way would increase our time to reach the goal. *Is it safe to turn left early?* Accelerating and turning left does not cause a collision. *Why not?* Because vehicle 3 stops. *What if vehicle 3 went straight?* We would slow down and give way. | |

in bold the explanations from CEMA. Second, to measure trust, we use the 9 trust scales proposed by Hoffman et al. [20] adapted to the AD domain. We use a between-subjects design: half the participants are shown the trust scales prior to ranking explanations,
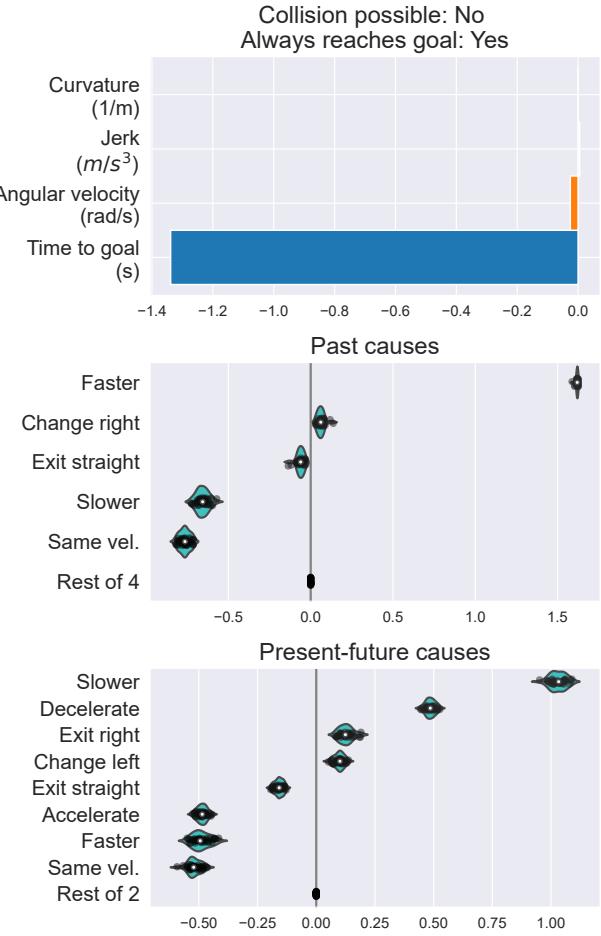


**Figure 4: [Top]** Signed differences between expected reward components correctly identify time-to-goal as the most significant teleological cause. **[Mid/Bot]** Feature importance attributions for the slice before and during/after the queried subsequence correctly rank mechanistic causes. Violin plots show 5-fold cross-validation repeated 7 times.

and the other half after having ranked explanations. We also asked participants about their driving experience and previous exposure to AVs using the SAE automation scale [36]. We hypothesize that **H1:** the explanations generated by CEMA are scored on average as highly as the human baseline explanations; **H2** participants who saw explanations from CEMA have on average higher levels of trust than those who have not. We analyze our data by fitting linear mixed-effects models for each hypothesis. We report the estimated means ($\hat{\beta}$) and standard errors ($\sigma$) for each variable and use the Wald test [43] to determine whether the effects of a variable are statistically significant on the outcome.

For **H1**, we found that CEMA's explanations were rated significantly higher when its explanations were not highlighted and were not significantly worse when they were highlighted. On average, explanation ratings ($\hat{\beta}_0$=3.31, $\sigma$=0.08) were marginally lower for human-written explanations ($\hat{\beta}$=−0.16, $\sigma$=0.08, $p$=0.21), and ratings
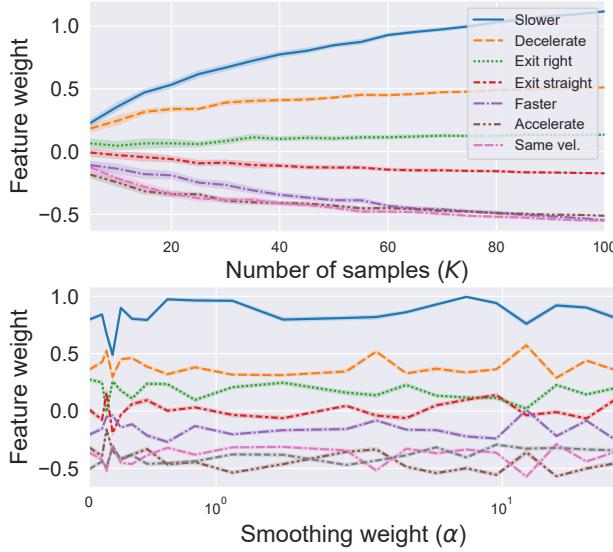
**Figure 5: Changes to causal attributions with [Top] different sample sizes and [Bot] different smoothing weights for present-future mechanistic causes in conversation S1-A. Shaded regions are bootstrapped 95% confidence intervals.**

were significantly lower for human explanations when CEMA's explanations were not highlighted to participants ($\hat{\beta}$=−0.22, $\sigma$=0.08, $p < 0.05$). Variations across scenarios were negligible ($SD$=0.07). We also found that people tend to rank CEMA's explanations higher when they had exposure to AVs previously ($\hat{\beta}$=0.1, $\sigma$=0.06, $p = 0.09$).

For **H2**, we found that, on average, participants' trust ratings ($\hat{\beta}_0$=1.53, $\sigma$=0.5) were significantly higher after seeing explanations ($\hat{\beta}$=0.11, $\sigma$=0.05, $p < 0.05$), which aligns with expectations from literature [30]. Participants' trust also increased significantly when they rated CEMA's explanations higher ($\hat{\beta}$=0.35, $\sigma$=0.15, $p \approx 0$) or when they had previous exposure to AVs ($\hat{\beta}$=0.33, $\sigma$=0.05, $p \approx 0$), but trust remained largely unchanged by human explanations ($\hat{\beta}$=0.12, $\sigma$=0.15, $p$=0.85). Trust ratings were not significantly affected by whether CEMA's explanations were highlighted ($\hat{\beta}$=0.02, $\sigma$=0.05, $p = 0.66$) and there were no significant interaction effects between the average ratings of explanations and highlighting ($\hat{\beta}$=−0.04, $\sigma$=0.04, $p = 0.36$). The estimated trust levels varied across the 9 trust scales ($SD$=0.53) but not the observed tendencies. Our results suggest that people who had some exposure to AVs or had a preference for CEMA's explanations were more likely to trust AVs in general, regardless of whether they knew which explanations came from CEMA. Taken together with the result for **H1**, this suggests that CEMA's explanations may be more effective at improving people's trust in AVs than non-expert human explanations.

## 7 DISCUSSION AND FUTURE WORK

Our primary goal with CEMA is to advance the field of social XAI applied to dynamic multi-agent systems. A crucial component of intelligible explanations is the use of semantically meaningful features [11]. Importantly, the challenge of designing useful features is not unique to CEMA but is a necessary step for any automated explanation generation system in social XAI. With CEMA, we assumed that there is a feature function $\phi$ which performs the translation from the raw representations of state and action spaces to the more abstract semantic feature space. This translation from state to feature space is domain-dependent and should be considered a crucial step during the deployment of social XAI systems. However, CEMA is feature-agnostic so that counterfactual causal selection does not depend on $\phi$ or the interpretations of features.

CEMA also does not rely on a fixed causal graph to model dynamic multi-agent systems. Instead, it assumes that there is a probabilistic model, such as a stochastic planner, trained joint policy, or autoregressive model trained on observational data, which can be used to forward simulate the state of the system. Based on the work of Quillien and Lucas [34] and the counterfactual model of causation [19, 24], CEMA can derive causes to an ego agent's actions in any system where such a model is obtainable. The assumption here is that these models cover alternatives that are grounded in factual observations with a non-zero probability, and any reasonably expressive algorithm would fulfill these criteria.

The user study suggests that people may prefer explanations generated by CEMA, however, trust levels are still low. This may be – as several participants indicated in their feedback – because people prefer to see agents act more conservatively, without exploiting potentially riskier but more efficient actions. Explanations that justify efficient but less safe decisions then have to overcome the inherent wariness of people, which was indeed high among participants, though it somewhat decreased after seeing explanations.

We designed CEMA to be used in conversations with users, but we did not focus on natural language processing in this work. For example, we assume that queries unambiguously describe the timing of actions – allowing us to focus on causal selection – but actual natural language queries are fuzzy and imprecise. By building modern NLP components, we can strengthen the social and conversational aspects of CEMA. Future work will involve the integration of language parsing [25] and dialogue systems [10] leveraging modern neural language models to deliver explanations.

Our implementation of CEMA for AD improves on existing social XAI methods for AD in several aspects. In contrast to Omeiza et al. [30], we avoid using a surrogate model and generate causal explanations that take the temporal nature of driving into account. Compared to Gyevnar et al. [14], CEMA supports multiple modes of explanations with both high-level and low-level features.

To conclude, our goal is to address some of the transparency-related social concerns of AI. CEMA fills a gap in social XAI by enabling causal explanation generation in dynamic sequential multi-agent systems. As we expect to see autonomous agents proliferate in everyday environments, social explanations will be crucial for building user trust and for the acceptance of new technologies.

# REFERENCES

[1] Stefano V. Albrecht, Cillian Brewitt, John Wilhelm, Balint Gyevnar, Francisco Eiras, Mihai Dobre, and Subramanian Ramamoorthy. 2021. Interpretable Goal-based Prediction and Planning for Autonomous Driving. In *IEEE International Conference on Robotics and Automation (ICRA)*.

[2] Stefano V. Albrecht and Peter Stone. 2018. Autonomous agents modelling other agents: A comprehensive survey and open problems. *Artificial Intelligence* 258 (2018), 66–95. https://doi.org/10.1016/j.artint.2018.01.002

[3] Sajid Ali, Tamer Abuhmed, Shaker El-Sappagh, Khan Muhammad, Jose M. Alonso-Moral, Roberto Confalonieri, Riccardo Guidotti, Javier Del Ser, Natalia Díaz-Rodríguez, and Francisco Herrera. 2023. Explainable Artificial Intelligence (XAI): What We Know and What Is Left to Attain Trustworthy Artificial Intelligence. *Information Fusion* (April 2023), 101805. https://doi.org/10.1016/j.inffus.2023.101805

[4] Peter C. Austin. 2011. An Introduction to Propensity Score Methods for Reducing the Effects of Confounding in Observational Studies. *Multivariate Behavioral Research* 46, 3 (May 2011), 399–424. https://doi.org/10.1080/00273171.2011.568786

[5] Claudine Badue, Rânik Guidolini, Raphael Vivacqua Carneiro, Pedro Azevedo, Vinicius B. Cardoso, Avelino Forechi, Luan Jesus, Rodrigo Berriel, Thiago M. Paixão, Filipe Mutz, Lucas de Paula Veronese, Thiago Oliveira-Santos, and Alberto F. De Souza. 2021. Self-Driving Cars: A Survey. *Expert Systems with Applications* 165 (March 2021), 113816. https://doi.org/10.1016/j.eswa.2020.113816

[6] Cillian Brewitt, Balint Gyevnar, Samuel Garcin, and Stefano V. Albrecht. 2021. GRIT: Fast, Interpretable, and Verifiable Goal Recognition with Learned Decision Trees for Autonomous Driving. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. 1023–1030. https://doi.org/10.1109/IROS51168.2021.9636279

[7] Cillian Brewitt, Massimiliano Tamborski, Cheng Wang, and Stefano V. Albrecht. 2023. Verifiable Goal Recognition for Autonomous Driving with Occlusions. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*.

[8] Nadia Burkart and Marco F. Huber. 2021. A Survey on the Explainability of Supervised Machine Learning. *Journal of Artificial Intelligence Research* 70 (May 2021), 245–317. https://doi.org/10.1613/jair.1.12228

[9] Tathagata Chakraborti, Sarath Sreedharan, and Subbarao Kambhampati. 2020. The Emerging Landscape of Explainable Automated Planning & Decision Making. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-20*, Christian Bessiere (Ed.). International Joint Conferences on Artificial Intelligence Organization, 4803–4811. https://doi.org/10.24963/ijcai.2020/669 Survey track.

[10] Hongshen Chen, Xiaorui Liu, Dawei Yin, and Jiliang Tang. 2017. A Survey on Dialogue Systems: Recent Advances and New Frontiers. *ACM SIGKDD Explorations Newsletter* 19, 2 (Nov. 2017), 25–35. https://doi.org/10.1145/3166054.3166058

[11] Richard Dazeley, Peter Vamplew, Cameron Foale, Charlotte Young, Sunil Aryal, and Francisco Cruz. 2021. Levels of Explainable Artificial Intelligence for Human-Aligned Conversational Explanations. *Artificial Intelligence* 299 (Oct. 2021), 103525. https://doi.org/10.1016/j.artint.2021.103525

[12] Albert Gatt and Ehud Reiter. 2009. SimpleNLG: A Realisation Engine for Practical Applications. In *Proceedings of the 12th European Workshop on Natural Language Generation (ENLG 2009)*. Association for Computational Linguistics, Athens, Greece, 90–93.

[13] Balint Gyevnar, Nick Ferguson, and Burkhard Schafer. 2023. Bridging the Transparency Gap: What Can Explainable AI Learn From the AI Act?. In *Proceedings of the 26th European Conference on Artificial Intelligence (ECAI 2023)*. IOS Press, 964–971. https://doi.org/10.3233/FAIA230367

[14] Balint Gyevnar, Massimiliano Tamborski, Cheng Wang, Christopher G. Lucas, Shay B. Cohen, and Stefano V. Albrecht. 2022. A Human-Centric Method for Generating Causal Explanations in Natural Language for Autonomous Vehicle Motion Planning. In *Workshop on Artificial Intelligence for Autonomous Driving*. International Joint Conference on Artificial Intelligence. https://doi.org/10.48550/arXiv.2206.08783 arXiv:2206.08783 [cs]

[15] Balint Gyevnar, Cheng Wang, Christopher G. Lucas, Shay B. Cohen, and Stefano V. Albrecht. 2024. HEADD: Human Explanations for Autonomous Driving Decisions. https://doi.org/10.7488/ds/7676

[16] R. J. Hankinson. 1998. *Cause and Explanation in Ancient Greek Thought*. Clarendon Press.

[17] Josiah P. Hanna, Arrasy Rahman, Elliot Fosong, Francisco Eiras, Mihai Dobre, John Redford, Subramanian Ramamoorthy, and Stefano V. Albrecht. 2021. Interpretable Goal Recognition in the Presence of Occluded Factors for Autonomous Vehicles. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*.

[18] Jacob Haspiel, Na Du, Jill Meyerson, Lionel P. Robert Jr., Dawn Tilbury, X. Jessie Yang, and Anuj K. Pradhan. 2018. Explanations and Expectations: Trust Building in Automated Vehicles. In *Companion of the 2018 ACM/IEEE International Conference on Human-Robot Interaction (HRI '18)*. Association for Computing Machinery, New York, NY, USA, 119–120. https://doi.org/10.1145/3173386.3177057

[19] Denis J. Hilton. 1988. Logic and Causal Attribution. In *Contemporary Science and Natural Explanation: Commonsense Conceptions of Causality*. New York University Press, New York, NY, US, 33–65.

[20] Robert R. Hoffman, Shane T. Mueller, Gary Klein, and Jordan Litman. 2019. Metrics for Explainable AI: Challenges and Prospects. *arXiv:1812.04608 [cs]* (Feb. 2019). arXiv:1812.04608 [cs]

[21] Renhao Huang, Hao Xue, Maurice Pagnucco, Flora Salim, and Yang Song. 2023. Multimodal Trajectory Prediction: A Survey. arXiv:2302.10463 [cs.RO]

[22] Davinder Kaur, Suleyman Uslu, Kaley J. Rittichier, and Arjan Durresi. 2022. Trustworthy Artificial Intelligence: A Review. *Comput. Surveys* 55, 2 (Jan. 2022), 39:1–39:38. https://doi.org/10.1145/3491209

[23] Jinkyu Kim, Anna Rohrbach, Trevor Darrell, John Canny, and Zeynep Akata. 2018. Textual Explanations for Self-Driving Vehicles. In *Computer Vision – ECCV 2018 (Lecture Notes in Computer Science)*, Vittorio Ferrari, Martial Hebert, Cristian Sminchisescu, and Yair Weiss (Eds.). Springer International Publishing, Cham, 577–593. https://doi.org/10.1007/978-3-030-01216-8_35

[24] David Lewis. 1973. Causation. *Journal of Philosophy* 70, 17 (1973), 556–567. https://doi.org/10.2307/2025310

[25] Jiaqi Li, Ming Liu, Bing Qin, and Ting Liu. 2022. A Survey of Discourse Parsing. *Frontiers of Computer Science* 16, 5 (Jan. 2022), 165329. https://doi.org/10.1007/s11704-021-0500-z

[26] Tania Lombrozo and Susan Carey. 2006. Functional Explanation and the Function of Explanation. *Cognition* 99, 2 (March 2006), 167–204. https://doi.org/10.1016/j.cognition.2004.12.009

[27] Prashan Madumal, Tim Miller, Liz Sonenberg, and Frank Vetere. 2020. Explainable Reinforcement Learning through a Causal Lens. *Proceedings of the AAAI Conference on Artificial Intelligence* 34, 03 (April 2020), 2493–2500. https://doi.org/10.1609/aaai.v34i03.5631

[28] Tim Miller. 2019. Explanation in Artificial Intelligence: Insights from the Social Sciences. *Artificial Intelligence* 267 (Feb. 2019), 1–38. https://doi.org/10.1016/j.artint.2018.07.007

[29] Samer B. Nashed, Saaduddin Mahmud, Claudia V. Goldman, and Shlomo Zilberstein. 2023. Causal Explanations for Sequential Decision Making Under Uncertainty. In *Proceedings of the 2023 International Conference on Autonomous Agents and Multiagent Systems (AAMAS '23)*. International Foundation for Autonomous Agents and Multiagent Systems, Richland, SC, 2307–2309.

[30] Daniel Omeiza, Helena Web, Marina Jirotka, and Lars Kunze. 2021. Towards Accountability: Providing Intelligible Explanations in Autonomous Driving. *Proceedings of the 32nd IEEE Intelligent Vehicles Symposium* (2021).

[31] OpenAI. 2022. ChatGPT: Optimizing Language Models for Dialogue. https://openai.com/blog/chatgpt/.

[32] Judea Pearl. 2009. *Causality* (second ed.). Cambridge University Press, Cambridge. https://doi.org/10.1017/CBO9780511803161

[33] Yunpeng Qing, Shunyu Liu, Jie Song, and Mingli Song. 2022. A Survey on Explainable Reinforcement Learning: Concepts, Algorithms, Challenges. arXiv:2211.06665 [cs]

[34] Tadeg Quillien and Christopher G. Lucas. 2023. Counterfactuals and the Logic of Causal Selection. *Psychological Review* Advance online publication (2023). https://doi.org/10.1037/rev0000428

[35] R Core Team. 2023. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. https://www.R-project.org

[36] SAE International. 2021. *Taxonomy and Definitions for Terms Related to Driving Automation Systems for On-Road Motor Vehicles*. Technical Report. United States.

[37] Wilko Schwarting, Javier Alonso-Mora, and Daniela Rus. 2018. Planning and Decision-Making for Autonomous Vehicles. *Annual Review of Control, Robotics, and Autonomous Systems* 1, 1 (2018), 187–210. https://doi.org/10.1146/annurev-control-060117-105157

[38] Andrew Silva, Matthew Gombolay, Taylor Killian, Ivan Jimenez, and Sung-Hyun Son. 2020. Optimization Methods for Interpretable Differentiable Decision Trees Applied to Reinforcement Learning. In *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*. PMLR, 1855–1865.

[39] Ilia Stepin, Jose M. Alonso, Alejandro Catala, and Martín Pereira-Fariña. 2021. A Survey of Contrastive and Counterfactual Explanation Generation Methods for Explainable Artificial Intelligence. *IEEE Access* 9 (2021), 11974–12001. https://doi.org/10.1109/ACCESS.2021.3051315

[40] Stratis Tsirtsis, Abir De, and Manuel Rodriguez. 2021. Counterfactual Explanations in Sequential Decision Making Under Uncertainty. In *Advances in Neural Information Processing Systems*, Vol. 34. Curran Associates, Inc., 30127–30139.

[41] Abhinav Verma, Vijayaraghavan Murali, Rishabh Singh, Pushmeet Kohli, and Swarat Chaudhuri. 2019. Programmatically Interpretable Reinforcement Learning. *arXiv:1804.02477 [cs, stat]* (April 2019). arXiv:1804.02477 [cs, stat]

[42] Sandra Wachter, Brent Mittelstadt, and Chris Russell. 2017. Counterfactual Explanations without Opening the Black Box: Automated Decisions and the GDPR. *Harvard Journal of Law & Technology (Harvard JOLT)* 31, 2 (2017), 841–888.

[43] Abraham Wald. 1943. Tests of Statistical Hypotheses Concerning Several Parameters When the Number of Observations is Large. *Trans. Amer. Math. Soc.* 54, 3 (1943), 426–482. http://www.jstor.org/stable/1990256

[44] Yiwen Zhang, Wenjia Wang, Xinyan Zhou, Qi Wang, and Xiaohua Sun. 2022. Tactical-Level Explanation Is Not Enough: Effect of Explaining AV's Lane-Changing Decisions on Drivers' Decision-Making, Trust, and Emotional Experience. *International Journal of Human–Computer Interaction* 0, 0 (Aug. 2022), 1–17. https://doi.org/10.1080/10447318.2022.2098965

## A IMPLEMENTATION DETAILS

### A.1 User Queries

As we our focus is on causal selection and NLP, we do not use NLU systems to parse users' questions. Instead, we define a standard query template and code queries as a structured JSON file. The template is given in Table 3 with each field and the corresponding type. The queries are given in terms of three main components: a query type (Type), a list of actions (Actions), and tense of the question (Tense); with an optional list of factual actions (Factuals). A list of actions is given as a combination of maneuvers and macro actions of IGP2. We identify for each timestep of a trajectory which maneuver and macro action are executed allowing for a timestep-by-timestep comparison to the queried list of actions.

The above information is used to infer the start $u$ and end timestep $v$ necessary to determine the state sequence $\hat{s}_{u:v}$. First, we use tense to filter out those timesteps of the observed trajectory $s^\varepsilon_{1:t}$ that do not match the tense of the question. For example, if the question is in future tense, then all past and present timesteps are removed. We then *group* consecutive timesteps by actions and filter out those action groups that do not match the user-given list of actions. If after this process, we still end up with multiple groups, then the first timestep of the group closest to the current timestep $t$ is picked as $u$ for the explanation. The end timestep $v$ is then given by the last timestep of the action group to which $u$ belongs.

**Table 3: The format template to encode queries. Types marked with a star are optional.**

| Field | Type | Explanation |
|---|---|---|
| Type | *str* | Type of the user's question. |
| VID (VID) | *int* | The VIDentifier for $\varepsilon$. |
| Tense | *str* | Grammatical tense of the user's question (past, present, future). |
| Actions | List[*str*] | A list of actions the user is interested in. |
| Query Time | *int* | Timestep of the user asking the question. |
| Action Time | *int** | Start time of queried actions. Only given for *what* query type. |
| Negated | *bool** | Whether the user's question is negated. |
| Factuals | List[*str*]* | The factual actions of $\varepsilon$. Only used if Actions is counterfactual. |

### A.2 Query Types

We also define three query types: *what*, *what if*, and *why* queries, which also affect the selected state sequence $\hat{s}_{u:v}$.

A *why*-query asks why a list of factual actions is executed. For example, *"Why did you change lanes to the left?"* It follows the process described above to determine $\hat{s}_{u:v}$. We concurrently generate both teleological and mechanistic explanations for this query type.

A *what if*-query asks about the actions of $\varepsilon$ had some other vehicle $i$ executed a different counterfactual list of actions. For example, *"What if vehicle 1 had stopped (instead of going straight)?"* Therefore, *what if*-queries encode a contrastive user question. We assume that the (counterfactual) list of actions directly replaces some factual list of actions that is also given for vehicle $i$ as part of the user's query. However, we cannot directly apply the process in Appendix A.1 to this list of actions to determine $\hat{s}_{u:v}$ as these actions have never been observed. Instead, we sample a set of alternative trajectories from the generative model $p(\hat{S}_{t+1:n}|s_{1:t})$ and find the counterfactual trajectory where vehicle $i$ was executing the counterfactual list of actions. We apply the process in Appendix A.1 to the selected counterfactual trajectory which filters out all action groups that do not correspond to the counterfactual list of actions. If multiple action groups remain, then we look at the additional factual list of actions from the query We use the observed trajectory and the factual list of actions to then find the factual action group of greatest overlap with any remaining counterfactual action groups. The values $u$ and $v$ are selected as the start and end timesteps of this overlapping region. This process guarantees that $u$ and $v$ happen at times when both the counterfactual and factual list of actions were executed in some possible worlds. We use the term *associative explanation* to refer to an explanation that describes the alternative actions of vehicle $\varepsilon$ had vehicle $i$ executed the counterfactual list of actions. We generate teleological and mechanistic explanations for why vehicle $\varepsilon$ would have executed those actions.

A *what*-query asks what vehicle $\varepsilon$ is doing at some given timestep. Thus, it always results in an associative explanation. Here, we do not have access to the list of actions since that is what we are trying to determine. Instead, we assume that the timestep $u$ of the start of the action (Action Time) is given.

Users' questions can be negated (Negated) sentences. The effect of negation is flipping the value of the boolean outcome variable $y$ in Algorithm 1, but it also affects how we determine $u$ and $v$. Negation turns a *why*-question into a *why not* question, as in *"Why did you not turn right?"* We treat this as a *what if* query asking about a counterfactual list of actions about vehicle $\varepsilon$ itself. For *what if* queries, negation means that a factual list of actions is already given, as in *"What if you hadn't stopped?"* Therefore, we can directly apply the process in Appendix A.1 to determine $u$ and $v$.

Finally, the user can query the future (e.g., future ego plan can be shown on a screen). For this, we concatenate the observed joint states $s_{1:t}$ with the *maximum a posteriori*-predictions of $p(\hat{S}_{t+1:n}|s_{1:t})$ giving $\hat{s}_{1:n}$. If the queried sequence is hypothetical, i.e., $\hat{s}_{u:v} \not\prec \hat{s}_{1:n}$, then we assume that we are given a corresponding factual subsequence to allow for the inference of the timings $u$ and $v$.

## A.3 Additive Smoothing

Since IGP2 can assign to some (reachable) goals and trajectories near-zero probabilities, we use additive smoothing with parameter $\alpha$ to make sure every trajectory can get sampled from $p(\hat{S}_{t+1:n}|s_{1:t})$. Given a discrete probability distribution $p_\theta : \Omega \mapsto [0, 1]$ parametrised by $\theta = [\theta_1, \ldots, \theta_d]$ where $\Omega$ is a finite non-empty set of events with size $d$, additive smoothing creates a new discrete probability distribution $p_\phi : \Omega \mapsto [0, 1]$ with new parameter vector $\phi = [\phi_1, \ldots, \phi_d]$ defined as:

$$\phi_i = \frac{\theta_i + \alpha}{1 + d\alpha}, \quad (i = 1, \ldots, d),$$

where we assumed that $\sum_i \theta_i = 1$.

## B EXPERIMENTAL SETUP

We implement CEMA using Python, building on the publicly available code repository of IGP2.[3] The source code for CEMA is available as supplementary material, is well-documented, and contains detailed instructions on how to reproduce our results. Our experiments were run on a modern Windows 10 PC with 32 GiB of RAM and a 12-core CPU. CEMA does not need a GPU.

Each scenario is defined using a semantic road layout in the ASAM OpenDrive 1.6 format, and using a configuration file that describes the behaviour and starting regions of agents on the road. Scenarios are executed in a simple, discrete-time simulation environment with an execution rate of 20 frames per second. For reproducing our results exactly, a random seed of 21 should be used whenever the algorithm uses randomness. Note, the results were confirmed with multiple seeds, we only fixed the seed for presentation in the paper.

Our scenarios are usually only a few seconds long (in simulation time), therefore, to avoid finding irrelevant causes and to speed up our implementation, we remove all states that are more than 5 seconds away from the current time step $t$ of the simulation. For all scenarios, we sample $K = 100$ counterfactual sequences with a smoothing weight $\alpha = 0.1$, and limits of $\tau_{min} = 2$ and $\tau_{max} = 5$ seconds. We picked $K$ so that it is large enough that a diverse range of trajectories is sampled considering that IGP2 predicts up to 3 distinct trajectories for each agent. The value for $\tau_{min}$ is the planning period of IGP2 (i.e., the number of seconds between two calls to IGP2), while $\tau_{max}$ is the maximum temporal distance as described above. We do not rely on any external datasets and no pre-processing steps are needed to run our code.

### B.1 Scaling with the Number of Agents

CEMA can scale to a large number of agents, which means that it can identify the relevant causes behind the actions of the ego agent, even if there are numerous irrelevant agents around, which did not have a causal effect on the actions of the ego. To show that this is true, we increase the number of agents in all four scenarios while making sure that the ego executes the same actions We then show that the same relevant causes, both teleological and mechanistic, are found by CEMA behind the ego's actions.

We alter the number of total agents in each scenario as follows: S1: 20 agents; S2: 16 agents; S3: 4 agents; S4: 16 agents. We only increase S3 to 5 agents, as the scenario map was too compact to

---

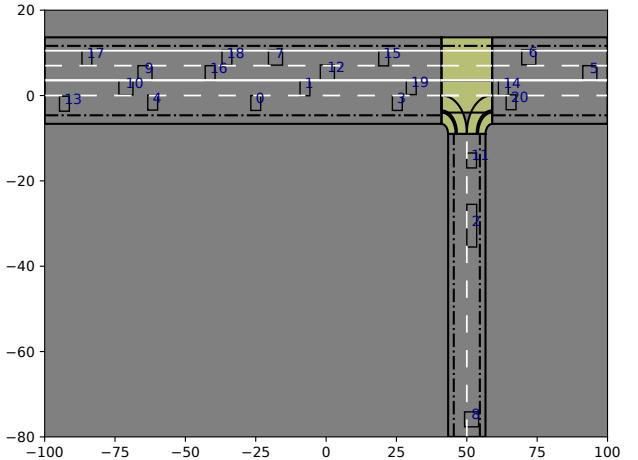[3]IGP2 available at https://uoe-agents.github.io/IGP2/.



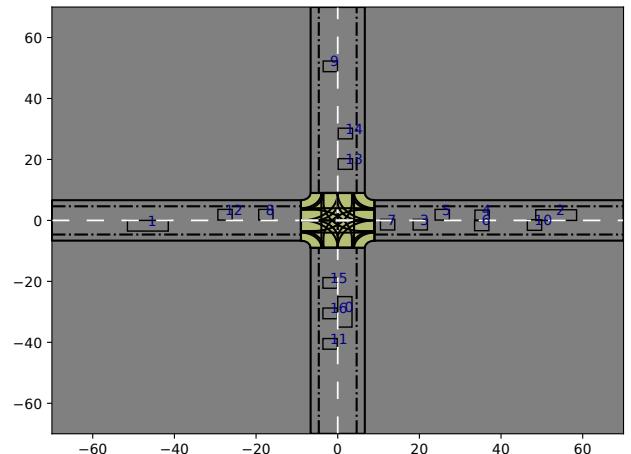Figure 6: Extended S1 with 20 spawn locations for agents.



Figure 7: S2 with 16 spawn locations for agents.

allow for more agents without a significant change to the behaviour of the ego. Figures 6 to 9 show the spawn regions of each agent, in which they are placed randomly.

### B.2 Size and Sampling Robustness

We test the size robustness of mechanistic explanations by sampling increasingly larger datasets using Algorithm 1 and for each dataset we rerun Algorithm 2 to determine causal attributions. We randomly sample $K \in \{5, 10, \ldots, 100\}$ different sample sizes. For a given $K$, we repeat the sampling process 50 times. We plot the mean and 95% bootstrapped confidence intervals of the causal attributions against $K$.

Using the fact that $\lim_{\alpha \to \infty} \phi_i = \frac{1}{d}$ for all $i$, we see that in the limit as $\alpha$ approaches infinity, $p_\phi$ approaches a uniform distribution with probability $\frac{1}{d}$ for all elements in $\Omega$. Thus we can interpolate between the original $p_\theta$ and an approximately uniform distribution by setting $\alpha$ to larger and larger values.
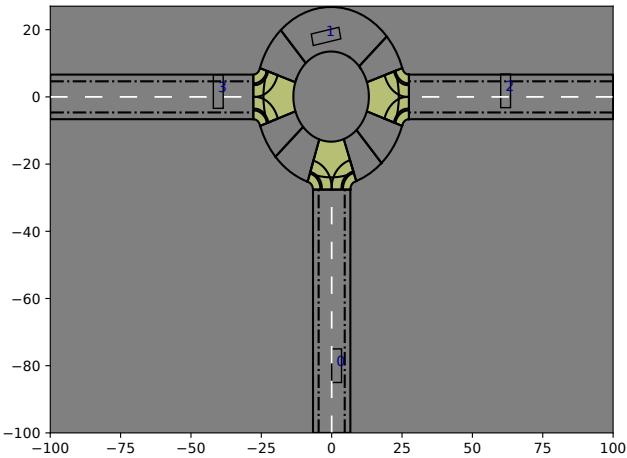
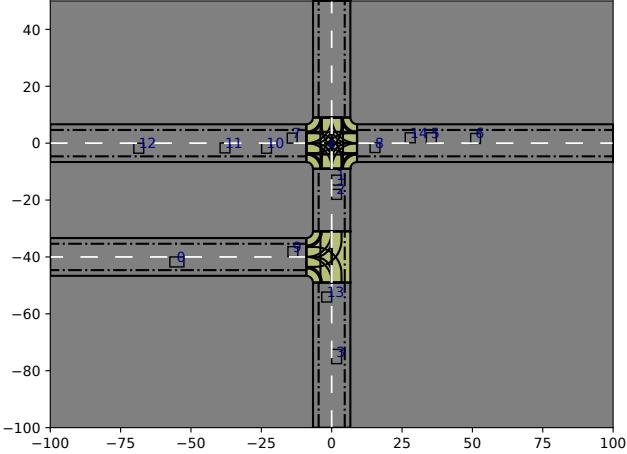Figure 8: S3 with 4 spawn locations for agents.



Figure 9: S4 with 16 spawn locations for agents.

**Table 4: The list of $\alpha$ values used for interpolation in the sampling robustness experiments.**

| |
|---|
| 0, 0.1, 0.14, 0.18, 0.25, 0.34, 0.45, 0.62, 0.83, 1.13, 1.53, 2.07, 2.8, 3.79, 5.13, 6.95, 9.41, 12.74, 17.25, 23.36, 31.62 |

Using the above, we perform our experiment for testing the robustness of our system against increasing uncertainty in external predictions of other vehicles' goals and trajectories. We define a list of twenty $\alpha$ values spaced on a logarithmic scale plus zero giving in total 21 distinct alphas (Table 4). For each $\alpha$ in this list, we sample a dataset of size $K = 50$ and determine causal attributions for mechanistic causes.

## C  COMPUTATIONAL EVALUATION

We present our full computational results for all queries in Figures 10 and 12 to 19. Figure 19 shows all *associative* (i.e., descriptive) queries. Figures 20 to 23 show the causal selection results for the scaling experiments described in Appendix B.1 Each figure gives natural language questions and the corresponding query content. We also give the full plot of causal attributions as well as the results for robustness experiments. In the plots of mechanistic causes, numbers in parentheses identify the ID of non-ego vehicles, if there is more than one non-ego. The code to reproduce these results is part of the supplementary material.

Details about the user study follows the results of this section.

**(a) Corresponding questions:** *"Why will you change lanes?"*; *"Why does it decrease the time to the goal?"*; *"Why will it be slower?"*

| Conv. | Type | VID | Tense | Actions | Query Time | Action Time | Negated | Factuals |
|-------|------|-----|-------|---------|-----------|-------------|---------|----------|
| S1-A | Why | 0 (AV) | Future | Change left | 40 | — | No | — |



Figure 10: Results for the first three questions in conversation S1-A.

**Figure 11: Corresponding question:** *"What if vehicle 1 hadn't changed lanes to the right before?"*

| Conv. | Type | VID | Tense | Actions | Query Time | Action Time | Negated | Factuals |
|-------|------|-----|-------|---------|-----------|-------------|---------|----------|
| S1-A | What if | 1 | Past | Change right | 75 | — | Yes | — |



(a) Reward difference

(c) Coefficient importance



**Figure 12: Results for the last question in conversation S1-A.**

| Conv. | Type | VID | Tense | Actions | Query Time | Action Time | Negated | Factuals |
|-------|------|-----|-------|---------|------------|-------------|---------|----------|
| S1-B | Why | 0 | Present | Go straight | 45 | — | Yes | Change left |



Figure 13: Results for the first three questions in conversation S1-B.

| Conv. | Type | VID | Tense | Actions | Query Time | Action Time | Negated | Factuals |
|-------|------|-----|-------|---------|------------|-------------|---------|----------|
| S2 | Why | 0 | Past | Give way & stop | 160 | — | Yes | Give way & accelerate |







Figure 14: Results for the first three questions in conversation S2.

| Conv. | Type | VID | Tense | Actions | Query Time | Action Time | Negated | Factuals |
|-------|------|-----|-------|---------|-----------|-------------|---------|----------|
| S2 | What if | 1 | Past | Go straight | 110 | — | No | Turn left |



Figure 15: Results for the last question in conversation S2.

| Conv. | Type | VID | Tense | Actions | Query Time | Action Time | Negated | Factuals |
|-------|------|-----|-------|---------|------------|-------------|---------|----------|
| S3 | Why | 0 | Present | Give way & stop | 90 | — | Yes | Give way & accelerate |



Figure 16: Results for the last two questions in conversation S3. There are multiple "change left" among the present-future causes as the non-ego could change lanes left at many points in the roundabout. CEMA differentiates between all of them, we only show it here this way for brevity.

| Conv. | Type | VID | Tense | Actions | Query Time | Action Time | Negated | Factuals |
|-------|------|-----|-------|---------|------------|-------------|---------|----------|
| S4 | Why | 0 | Present | Give way & stop | 140 | — | Yes | Turn left & accelerate |



Figure 17: Results for the first two questions in conversation S4.

| Conv. | Type | VID | Tense | Actions | Query Time | Action Time | Negated | Factuals |
|-------|------|-----|-------|---------|------------|-------------|---------|----------|
| S4 | What if | 3 | Future | Go straight | 100 | — | No | Stop |

Collision possible: No
Always reaches goal: Yes



No past causes because this is a what-if query.

Present-future causes

(a) Reward difference

(c) Coefficient importance



Present-future causes



**Figure 18: Results for the last question in conversation S4.**

**Figure 19: All queries that provide an *associative* explanation; across all conversations.**

| Conv. | Type | VID | Tense | Actions | Query Time | Action Time | Negated | Factuals |
|-------|------|-----|-------|---------|------------|-------------|---------|----------|
| S1-A | What if | 1 | Past | Change right | 75 | — | Yes | — |
| S1-B | What | 0 | — | — | 45 | 70 | No | — |
| S2 | What if | 0 | Past | Go straight | 110 | — | No | Turn left |
| S3 | What | 0 | — | — | 105 | 80 | No | — |
| S4 | What if | 1 | Future | Go straight | 100 | — | No | Stop |

**Figure 20: Top 10 causal importances in the extended scenario 1 with 20 agents.**



**Figure 21: Top 10 causal importances in the extended scenario 2 with 16 agents.**



**Figure 22: Top 10 causal importances in the extended scenario 3 with 4 agents.**



**Figure 23: Top 10 causal importances in the extended scenario 4 with 16 agents.**

# D USER STUDY

## D.1 Survey Methodology

We performed two surveys as part of the user study. The first elicited high-quality explanations from participants about the behaviour of the autonomous vehicle as shown in pre-recorded video clips. The second survey was used to compare the explanations of CEMA against this collected baseline dataset of human-written explanations. We include the dataset of high-quality human explanations, responses to the second survey, and code to reproduce our statistical analysis of the second survey as part of the supplementary materials.

*D.1.1 Participants.* We used the crowd-sourcing platform Prolific to recruit participants. We recruited from the USA, as the video recordings were in right-handed traffic, and filtered for participants whose first language was English. Participants were paid an average of £11/h and were shown our ethics approval and a consent form which they had to accept before being allowed to fill out the surveys.

Regarding the first survey, 54 participants filled it out with a median duration of completion of 25 minutes and 37 seconds. The sex distribution was 25 males and 29 females. The participants' age ranged between 19 to 73 years, with a median of 36 years. No participants were excluded for failing attention checks in the first survey.

Regarding the second survey, 200 participants filled it out with a median duration of completion of 9 minutes and 34 seconds. Sexes were distributed as 99 males and 101 females. The participants' age ranged between 19 to 77 years with a median of 35 years. One participant was excluded for failing attention checks in the second survey.

*D.1.2 Survey 1: Design and Procedure.* In the first survey, people were shown 7 scenarios out of a collection of 14 including the four scenarios used in the evaluation of CEMA. For each scenario, people were told what the goal of the ego agent is, and then they were shown a short (about 5 to 15-second-long) video of the scenario. The videos were top-down animations recorded in the software Road-Runner 2023a by MathWorks. We include the videos for the four scenarios used in our evaluation in the main paper. After watching the videos, participants were asked to answer the following four questions:

(1) Describe the actions of the blue car, self-driving car.
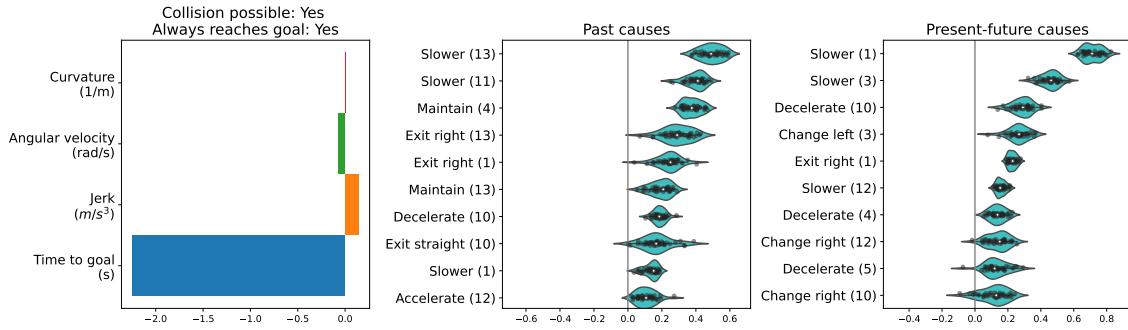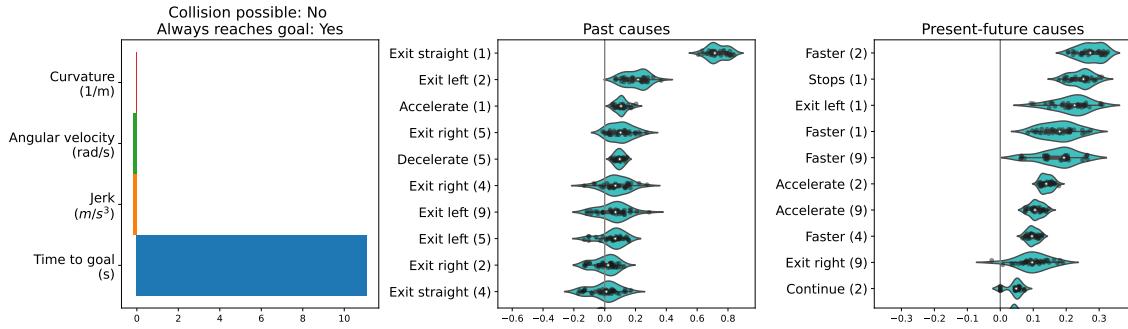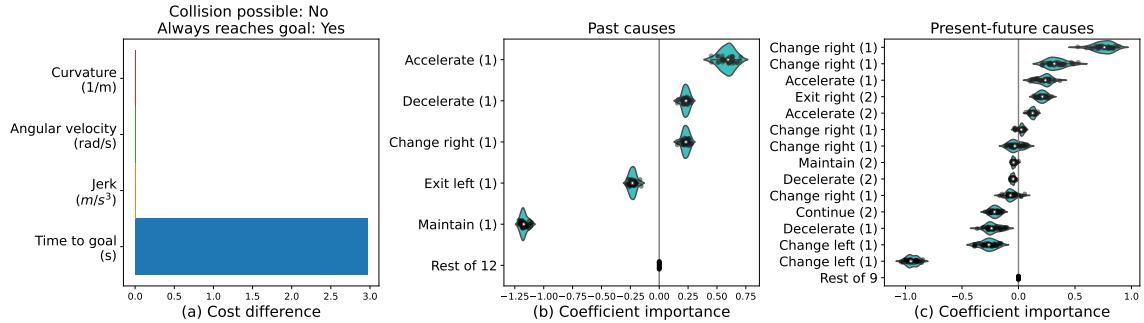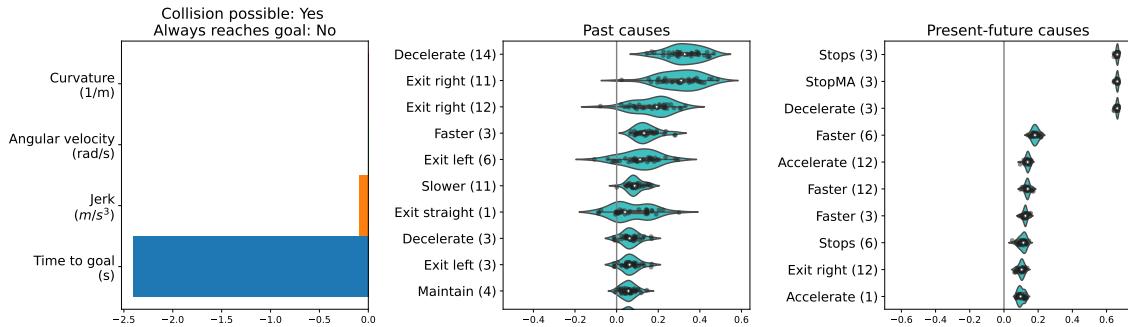(2) Explain why the blue car, self-driving car took these actions over different actions to reach its goal.
(3) Explain how the blue car, self-driving car was influenced in the scenario to take these actions.
(4) Describe changes to the scenario so that the blue car, self-driving car takes different actions? *(The new actions need not be the best actions in the original scenario)*

With 50% chance, participants were shown either "car" or "self-driving" car. We did this in order to get wide coverage of explanations, in case people were *a priori* biased against autonomous vehicles. Participants had to answer all four questions in their own words, entering their responses into unconstrained free-text boxes. We also introduced two attention checks at the end of the page of random scenarios #3 and #6.

*D.1.3 Creating the Baseline Dataset.* Having collected explanations in Survey 1, we needed to assemble a useful and high-quality baseline dataset for the second survey. As the second survey is only concerned with our four main scenarios, we excluded all explanations from Survey 1 that were not about scenarios S1-4. This left a total of 408 explanations. From this, we further removed another 26 explanations which were vacuous, that is, they did not contain any useful information about the scenario, or they expressed highly personal or obviously incorrect opinions.

To select explanations for the baseline dataset for Survey 2, we set up an annotation scheme with the following elements:

- Over quality: Scale on 1 to 5. The overall quality of an explanation as determined by its relevancy to the scenario, its linguistic correctness, and subjective clarity;
- Complexity: Scale on 1 to 5. The complexity of the explanation as determined by its length and the number of causal relationships mentioned in it;
- Type: Whether the explanation is a teleological (1), mechanistic (2), or associative (i.e., descriptive) (3) explanation, or some combination of all three;
- Counterfactual: Boolean. Whether the explanation is contrastive or not.

We annotated all remaining 382 explanations with this scheme and randomly picked 25 explanations for each of the four scenarios S1-4. We also corrected any obvious typos in each explanation. Our selection criteria included only the highest quality explanations (level 4 or 5) and a wide range of combinations of complexities, types, and counterfactuals was covered as well. This process in the end resulted in a selected dataset of high-quality baseline human-written explanations with a total size of 100 explanations.

*D.1.4 Survey 2: Design and Procedure.* We used Survey 2 to measure the effects of explanations on trust and to understand what kind of explanations people preferred: the ones generated by CEMA or the human-written baseline explanations. We randomised several parts of this survey, so for later reference, we will give names to random variables starting with an uppercase letter (shown in italics in parentheses), for example, *Type*.

Before the start of the survey (after giving consent), we asked participants to answer three questions to measure their previous exposure to autonomous driving systems (*AVExperience*), as classified by the SAE automation scale [36]. The question were:

- "Have you ever used a driver-assistance system? Examples of driver-assistance systems are adaptive cruise control, lane keeping assistance, and parking assistance." (SAE L1-2);
- "Have you ever used a partially autonomous self-driving car? A partially autonomous self-driving car is one which can drive on its own, but still requires an attentive human driver. An example is the Tesla Autopilot." (SAE L3);
- "Have you ever used a fully autonomous self-driving car? A fully autonomous self-driving car is one which does not require a human driver. An example is the Waymo One ride-hailing service." (SAE L4).

The core of Survey 2 was the explanation rating scales. For each participant, we randomly and with equal chance picked two of the four scenarios (*Scenario*), as well as 10 explanations to be ranked.

The explanations were picked from a set of 30 total explanations per scenario consisting of 5 CEMA-generated explanations the 25 selected explanations from Survey 1 (*Type*). For each scenario, similarly to Survey 1, we first showed each participant the video recording of the scenario, though this time we did not tell them what the goal of the ego vehicle was to avoid biasing people's ratings. We then asked participants to rate each explanation on a 5-point Likert scale based on how well each statement explained the actions of the ego vehicle (*MeanRating*). With 50% chance, we also highlighted to participants in bold typeface the explanations that were generated by CEMA (*Highlight*), and made sure that they were informed that these explanations were generated by the ego agent itself.

To measure changes to trust levels, we created the second survey with a between-subjects design, where the "treatment" was showing the participants the explanations and having them engage with the explanations through their ratings. Therefore, we assigned each participant randomly with a 50% chance into one of two groups: either answering the trust questions before seeing the explanations or after (*PreExplanation*). To actually measure people's trust (*TrustRating*) in AVs, we used the trust scales recommended by Hoffman et al. [20] adapted to the AD domain. People were asked the following nine questions in some random order (*TrustQuestion*) and asked to rate them on a 5-point Likert scale:

(1) "I am confident in self-driving cars. I feel that they work well" (confidence in AVs).
(2) "The decisions of self-driving cars are predictable." (AV decisions are predictable)
(3) "Self-driving cars are reliable. I can count on them to be correct at all times." (AVs are reliable)
(4) "I feel safe relying on a self-driving car. It will make the right decisions." (AVs are always correct)
(5) "Self-driving cars are efficient, in that they make decisions very quickly." (AVs are efficient)
(6) "I am wary of self-driving cars." (Wary of AVs)
(7) "Self-driving cars make better decisions than novice drivers." (AVs better than novice)
(8) "Self-driving cars make better decisions than experienced drivers." (AVs better than expert)
(9) "I would like to use self-driving cars for decision making." (Willing to AVs)

*D.1.5 Post-Survey Questions.* At the end of both surveys, participants were asked to fill out a brief survey about their driving experience about whether they hold a valid driver's license (*License*), how many years of driving experience they have (*Experience*), how many days a week they drive (*Frequency*), and how many miles on average they drive in a year (*Distance*). We also asked them voluntary demographic questions about their age range (*Age*), gender (*Gender*), and education level (*Education*). Finally, participants had the opportunity at the end of each survey to give any manner of feedback they thought worthy of mentioning.

## D.2   Analysis

To understand the effects of explanations and the ratings of different types of explanations, we conducted a detailed statistical analysis

of the survey data. For our analysis, we used the R programming language [35].

*D.2.1   Pre-processing.* We pre-processed the raw survey data to make it more amenable to analysis with R. This involved calculating the mean ratings for each participant grouped by the type of explanation (CEMA or Human) and the scenario (Scenario). It also involved collating the trust responses across participants into a single column grouped by the trust scale (TrustQuestion) and turning the data into long format. Finally, in addition to the raw explanation ratings, we also included fields for the mean explanation ratings of participants grouped by their Type (*MeanCEMARating* and *MeanHumanRating*).

*D.2.2   Modelling.* We used Gaussian-family mixed-effects linear models to understand the relationships between the various variables.

To understand how participants' ratings of explanations varied with different variables (corresponding to testing **H1**), we fitted a model to our data predicting the mean explanation ratings (MeanRating) of participants. We included the crossing between HighlightAV and Type, and AVExperience as independent variables encoding them using dummy coding. We set the reference levels: HighlightAV=TRUE, Type=CEMA, AVExperience=0 (no experience). The model also included varying intercepts for each scenario.

To understand how participants' trust levels changed after seeing explanations (corresponding to testing **H2**), we predicted the mean trust ratings (TrustRating) of participants. We included as factors, PreExplanation, HighlightAV, AVExperience, and the crossing between MeanCEMARating and MeanHumanRating. We encoded the categorical factors with dummy coding, setting the reference levels as follows: PreExplanation=TRUE, HighlightAV=TRUE, AVExperience=0. We also included varying intercepts grouped by TrustQuestion.

We have also tested more complex models with more variables, however, there was no significant change in results by including these, and the model would often not converge, so we opted to use the least complex but most expressive model, which we have reported here.

## D.3   Results

The following section presents a detailed view of the quantitative results of our analysis of Survey 2.

*D.3.1   Summary Statistics.* We present in Tables 5 and 6 the summary statistics for the trust levels. In Tables 7 to 10, we give the summary statistics of the ratings of all 30 questions for each scenario. The questions in bold were generated by CEMA. We also plot the distributions of trust before and after seeing explanations in Figures 24 and 25.

**Table 5: Trust levels before seeing explanations (PreExplanation = TRUE). Trust scale corresponding to the enumeration in Appendix D.1.4.**

| trust scale | count | mean | std | min | 25% | 50% | 75% | max |
|---|---|---|---|---|---|---|---|---|
| 1 | 101.0 | 2.653465 | 1.117458 | 1.0 | 2.00 | 2.0 | 4.00 | 5.0 |
| 2 | 101.0 | 3.099010 | 1.081711 | 1.0 | 2.00 | 3.0 | 4.00 | 5.0 |
| 3 | 101.0 | 2.297030 | 1.072796 | 1.0 | 1.00 | 2.0 | 3.00 | 5.0 |
| 4 | 101.0 | 2.425743 | 1.116661 | 1.0 | 2.00 | 2.0 | 3.00 | 5.0 |
| 5 | 101.0 | 3.297030 | 1.091280 | 1.0 | 3.00 | 3.0 | 4.00 | 5.0 |
| 6 | 101.0 | 3.881188 | 1.022616 | 1.0 | 3.00 | 4.0 | 5.00 | 5.0 |
| 7 | 101.0 | 3.217822 | 1.109991 | 1.0 | 3.00 | 3.0 | 4.00 | 5.0 |
| 8 | 101.0 | 2.376238 | 1.164916 | 1.0 | 1.00 | 2.0 | 3.00 | 5.0 |
| 9 | 101.0 | 2.415842 | 1.168480 | 1.0 | 2.00 | 2.0 | 3.00 | 5.0 |

**Table 6: Trust levels after seeing explanations (PreExplanation = FALSE). Trust scale corresponding to the enumeration in Appendix D.1.4.**

| trust scale | count | mean | std | min | 25% | 50% | 75% | max |
|---|---|---|---|---|---|---|---|---|
| 1 | 99.0 | 2.696970 | 1.092222 | 1.0 | 2.00 | 2.0 | 3.00 | 5.0 |
| 2 | 99.0 | 3.333333 | 1.133893 | 1.0 | 3.00 | 4.0 | 4.00 | 5.0 |
| 3 | 99.0 | 2.323232 | 1.095709 | 1.0 | 1.50 | 2.0 | 3.00 | 5.0 |
| 4 | 99.0 | 2.575758 | 1.107405 | 1.0 | 2.00 | 2.0 | 4.00 | 5.0 |
| 5 | 99.0 | 3.434343 | 0.949346 | 1.0 | 3.00 | 4.0 | 4.00 | 5.0 |
| 6 | 99.0 | 3.838384 | 1.037138 | 1.0 | 3.00 | 4.0 | 5.00 | 5.0 |
| 8 | 99.0 | 3.303030 | 1.063826 | 1.0 | 3.00 | 3.0 | 4.00 | 5.0 |
| 9 | 99.0 | 2.545455 | 1.032855 | 1.0 | 2.00 | 3.0 | 3.00 | 5.0 |
| 10 | 99.0 | 2.595960 | 1.105821 | 1.0 | 2.00 | 2.0 | 3.50 | 5.0 |



Figure 24: Mean ratings of participants on the 9 trust scales before and after having seen explanations with standard error.



Figure 25: Ratings of CEMA's and human explanations with bootstrapped 95% confidence interval (CI) grouped by whether CEMA's explanations were highlighted. The effect of highlighting is significant under our mixed-effects model ($p < 0.05$).

Table 7: Base statistics of explanation ratings for the selected 30 questions in Scenario 1.

| question | count | mean | std | min | 25% | 50% | 75% | max |
|---|---|---|---|---|---|---|---|---|
| 1 | 100.0 | 2.640000 | 1.267464 | 1.0 | 1.0 | 3.0 | 3.00 | 5.0 |
| 2 | 26.0 | 3.461538 | 1.103839 | 1.0 | 3.0 | 4.0 | 4.00 | 5.0 |
| 3 | 100.0 | 3.730000 | 1.081105 | 1.0 | 3.0 | 4.0 | 5.00 | 5.0 |
| 4 | 26.0 | 2.807692 | 1.059027 | 1.0 | 2.0 | 3.0 | 3.75 | 5.0 |
| 5 | 100.0 | 4.100000 | 1.010051 | 2.0 | 3.0 | 4.0 | 5.00 | 5.0 |
| 6 | 26.0 | 3.384615 | 0.803837 | 2.0 | 3.0 | 3.0 | 4.00 | 5.0 |
| 7 | 23.0 | 3.217391 | 1.241572 | 1.0 | 2.5 | 3.0 | 4.00 | 5.0 |
| 8 | 26.0 | 3.153846 | 1.255143 | 1.0 | 2.0 | 3.0 | 4.00 | 5.0 |
| 9 | 25.0 | 3.400000 | 1.118034 | 1.0 | 3.0 | 3.0 | 4.00 | 5.0 |
| 10 | 26.0 | 2.653846 | 1.198075 | 1.0 | 2.0 | 3.0 | 3.00 | 5.0 |
| 11 | 27.0 | 3.259259 | 1.163304 | 1.0 | 2.0 | 3.0 | 4.00 | 5.0 |
| 12 | 27.0 | 3.777778 | 1.012739 | 2.0 | 3.0 | 4.0 | 4.50 | 5.0 |
| 13 | 24.0 | 2.541667 | 1.062367 | 1.0 | 2.0 | 2.0 | 3.25 | 4.0 |
| 14 | 24.0 | 2.958333 | 1.334465 | 1.0 | 2.0 | 3.0 | 4.00 | 5.0 |
| 15 | 26.0 | 3.076923 | 1.262476 | 1.0 | 2.0 | 3.0 | 4.00 | 5.0 |
| 16 | 27.0 | 3.518519 | 1.087353 | 1.0 | 3.0 | 4.0 | 4.00 | 5.0 |
| 17 | 28.0 | 3.428571 | 0.997351 | 1.0 | 3.0 | 3.0 | 4.00 | 5.0 |
| 18 | 26.0 | 2.923077 | 1.293772 | 1.0 | 2.0 | 3.0 | 4.00 | 5.0 |
| 19 | 27.0 | 3.333333 | 1.176697 | 1.0 | 2.0 | 3.0 | 4.00 | 5.0 |
| 20 | 24.0 | 2.958333 | 1.122078 | 1.0 | 2.0 | 3.0 | 4.00 | 5.0 |
| 21 | 27.0 | 2.851852 | 0.988538 | 1.0 | 2.0 | 3.0 | 3.00 | 5.0 |
| 22 | 26.0 | 3.384615 | 1.202561 | 1.0 | 3.0 | 3.5 | 4.00 | 5.0 |
| 23 | 27.0 | 3.962963 | 0.939782 | 2.0 | 3.0 | 4.0 | 5.00 | 5.0 |
| 24 | 27.0 | 4.111111 | 1.219500 | 1.0 | 3.5 | 5.0 | 5.00 | 5.0 |
| 25 | 26.0 | 3.153846 | 1.120439 | 1.0 | 3.0 | 3.0 | 4.00 | 5.0 |
| 26 | 27.0 | 3.740741 | 0.813000 | 2.0 | 3.0 | 4.0 | 4.00 | 5.0 |
| 27 | 25.0 | 2.800000 | 1.258306 | 1.0 | 2.0 | 3.0 | 3.00 | 5.0 |
| 28 | 27.0 | 3.629630 | 1.005682 | 2.0 | 3.0 | 4.0 | 4.00 | 5.0 |
| 29 | 25.0 | 3.560000 | 0.916515 | 1.0 | 3.0 | 4.0 | 4.00 | 5.0 |
| 30 | 25.0 | 1.800000 | 0.912871 | 1.0 | 1.0 | 2.0 | 2.00 | 4.0 |

Table 8: Base statistics of explanation ratings for the selected 30 questions in Scenario 2.

| question | count | mean | std | min | 25% | 50% | 75% | max |
|---|---|---|---|---|---|---|---|---|
| 1 | 101.0 | 2.861386 | 1.191887 | 1.0 | 2.00 | 3.0 | 4.00 | 5.0 |
| 2 | 27.0 | 3.555556 | 1.250641 | 1.0 | 3.00 | 4.0 | 4.00 | 5.0 |
| 3 | 101.0 | 3.316832 | 1.264363 | 1.0 | 2.00 | 4.0 | 4.00 | 5.0 |
| 4 | 27.0 | 2.740741 | 1.163304 | 1.0 | 2.00 | 3.0 | 4.00 | 5.0 |
| 5 | 101.0 | 3.386139 | 1.288179 | 1.0 | 2.00 | 4.0 | 4.00 | 5.0 |
| 6 | 24.0 | 1.750000 | 1.032094 | 1.0 | 1.00 | 1.0 | 2.25 | 4.0 |
| 7 | 25.0 | 2.640000 | 1.319091 | 1.0 | 1.00 | 3.0 | 3.00 | 5.0 |
| 8 | 25.0 | 3.440000 | 1.157584 | 1.0 | 3.00 | 4.0 | 4.00 | 5.0 |
| 9 | 26.0 | 3.769231 | 1.274604 | 1.0 | 3.00 | 4.0 | 5.00 | 5.0 |
| 10 | 27.0 | 3.740741 | 1.195909 | 1.0 | 3.50 | 4.0 | 4.50 | 5.0 |
| 11 | 26.0 | 3.076923 | 1.293772 | 1.0 | 2.00 | 3.0 | 4.00 | 5.0 |
| 12 | 25.0 | 3.160000 | 1.312758 | 1.0 | 2.00 | 3.0 | 4.00 | 5.0 |
| 13 | 27.0 | 3.222222 | 1.368136 | 1.0 | 2.50 | 3.0 | 4.00 | 5.0 |
| 14 | 26.0 | 2.884615 | 1.423430 | 1.0 | 2.00 | 3.0 | 4.00 | 5.0 |
| 15 | 26.0 | 3.846154 | 0.967153 | 1.0 | 3.25 | 4.0 | 4.00 | 5.0 |
| 16 | 27.0 | 2.851852 | 1.199478 | 1.0 | 2.00 | 3.0 | 4.00 | 5.0 |
| 17 | 26.0 | 3.153846 | 1.461296 | 1.0 | 2.00 | 3.5 | 4.00 | 5.0 |
| 18 | 27.0 | 3.518519 | 1.369176 | 1.0 | 2.00 | 4.0 | 4.50 | 5.0 |
| 19 | 26.0 | 1.538462 | 0.859338 | 1.0 | 1.00 | 1.0 | 2.00 | 3.0 |
| 20 | 25.0 | 2.960000 | 1.513275 | 1.0 | 2.00 | 3.0 | 5.00 | 5.0 |
| 21 | 26.0 | 2.923077 | 1.262476 | 1.0 | 2.00 | 3.0 | 4.00 | 5.0 |
| 22 | 26.0 | 3.269231 | 1.372813 | 1.0 | 2.00 | 3.5 | 4.00 | 5.0 |
| 23 | 25.0 | 3.080000 | 1.222020 | 1.0 | 2.00 | 3.0 | 4.00 | 5.0 |
| 24 | 27.0 | 2.555556 | 1.527525 | 1.0 | 1.00 | 2.0 | 4.00 | 5.0 |
| 25 | 26.0 | 3.500000 | 0.948683 | 2.0 | 3.00 | 3.5 | 4.00 | 5.0 |
| 26 | 27.0 | 2.925926 | 1.327981 | 1.0 | 2.00 | 3.0 | 4.00 | 5.0 |
| 27 | 27.0 | 3.740741 | 1.403090 | 1.0 | 3.00 | 4.0 | 5.00 | 5.0 |
| 28 | 27.0 | 3.222222 | 1.086042 | 1.0 | 2.00 | 4.0 | 4.00 | 5.0 |
| 29 | 27.0 | 2.925926 | 1.412198 | 1.0 | 2.00 | 3.0 | 4.00 | 5.0 |
| 30 | 27.0 | 3.111111 | 1.476309 | 1.0 | 2.00 | 3.0 | 4.00 | 5.0 |

**Table 9: Base statistics of explanation ratings for the selected 30 questions in Scenario 3.**

| question | count | mean | std | min | 25% | 50% | 75% | max |
|---|---|---|---|---|---|---|---|---|
| 1 | 98.0 | 2.857143 | 1.192762 | 1.0 | 2.00 | 3.0 | 4.00 | 5.0 |
| 2 | 26.0 | 3.307692 | 1.123182 | 1.0 | 3.00 | 3.0 | 4.00 | 5.0 |
| 3 | 98.0 | 3.693878 | 1.106757 | 1.0 | 3.00 | 4.0 | 5.00 | 5.0 |
| 4 | 25.0 | 2.600000 | 1.258306 | 1.0 | 2.00 | 3.0 | 3.00 | 5.0 |
| 5 | 98.0 | 3.163265 | 1.289880 | 1.0 | 2.00 | 3.0 | 4.00 | 5.0 |
| 6 | 23.0 | 2.478261 | 1.591728 | 1.0 | 1.00 | 3.0 | 3.50 | 5.0 |
| 7 | 26.0 | 3.346154 | 1.263085 | 1.0 | 2.00 | 4.0 | 4.00 | 5.0 |
| 8 | 26.0 | 3.884615 | 1.275207 | 1.0 | 3.00 | 4.0 | 5.00 | 5.0 |
| 9 | 26.0 | 3.653846 | 1.354764 | 1.0 | 3.00 | 4.0 | 5.00 | 5.0 |
| 10 | 25.0 | 3.320000 | 1.314027 | 1.0 | 2.00 | 3.0 | 4.00 | 5.0 |
| 11 | 25.0 | 3.360000 | 1.254326 | 1.0 | 3.00 | 3.0 | 5.00 | 5.0 |
| 12 | 24.0 | 3.458333 | 1.215092 | 1.0 | 3.00 | 4.0 | 4.00 | 5.0 |
| 13 | 25.0 | 3.160000 | 1.247664 | 1.0 | 2.00 | 3.0 | 4.00 | 5.0 |
| 14 | 26.0 | 2.192308 | 1.414757 | 1.0 | 1.00 | 2.0 | 3.00 | 5.0 |
| 15 | 26.0 | 3.692308 | 1.378963 | 1.0 | 3.00 | 4.0 | 5.00 | 5.0 |
| 16 | 26.0 | 2.730769 | 1.457606 | 1.0 | 1.00 | 3.0 | 4.00 | 5.0 |
| 17 | 26.0 | 2.807692 | 1.497177 | 1.0 | 1.25 | 3.0 | 4.00 | 5.0 |
| 18 | 26.0 | 3.692308 | 1.225373 | 1.0 | 3.00 | 4.0 | 5.00 | 5.0 |
| 19 | 24.0 | 2.333333 | 1.403928 | 1.0 | 1.00 | 2.0 | 3.25 | 5.0 |
| 20 | 23.0 | 3.478261 | 1.441892 | 1.0 | 2.50 | 4.0 | 5.00 | 5.0 |
| 21 | 26.0 | 2.692308 | 1.257592 | 1.0 | 2.00 | 3.0 | 3.00 | 5.0 |
| 22 | 27.0 | 3.925926 | 1.106829 | 1.0 | 3.00 | 4.0 | 5.00 | 5.0 |
| 23 | 25.0 | 3.760000 | 1.267544 | 1.0 | 3.00 | 4.0 | 5.00 | 5.0 |
| 24 | 25.0 | 3.040000 | 1.171893 | 1.0 | 2.00 | 3.0 | 4.00 | 5.0 |
| 25 | 27.0 | 3.814815 | 1.075498 | 1.0 | 3.00 | 4.0 | 5.00 | 5.0 |
| 26 | 26.0 | 2.730769 | 1.218448 | 1.0 | 2.00 | 3.0 | 3.75 | 5.0 |
| 27 | 26.0 | 2.538462 | 1.475961 | 1.0 | 1.00 | 2.0 | 4.00 | 5.0 |
| 28 | 25.0 | 2.520000 | 1.326650 | 1.0 | 1.00 | 2.0 | 3.00 | 5.0 |
| 29 | 26.0 | 2.115385 | 1.107318 | 1.0 | 1.00 | 2.0 | 3.00 | 4.0 |
| 30 | 25.0 | 3.680000 | 1.069268 | 2.0 | 3.00 | 4.0 | 5.00 | 5.0 |

Table 10: Base statistics of explanation ratings for the selected 30 questions in Scenario 4.

| question | count | mean | std | min | 25% | 50% | 75% | max |
|----------|-------|----------|----------|-----|------|-----|------|-----|
| 1 | 101.0 | 2.871287 | 1.270145 | 1.0 | 2.00 | 3.0 | 4.00 | 5.0 |
| 2 | 27.0 | 3.629630 | 1.079464 | 1.0 | 3.00 | 4.0 | 4.00 | 5.0 |
| 3 | 25.0 | 3.560000 | 1.083205 | 2.0 | 3.00 | 3.0 | 5.00 | 5.0 |
| 4 | 101.0 | 2.900990 | 1.360183 | 1.0 | 2.00 | 3.0 | 4.00 | 5.0 |
| 5 | 101.0 | 3.722772 | 1.078136 | 1.0 | 3.00 | 4.0 | 5.00 | 5.0 |
| 6 | 27.0 | 2.777778 | 1.154701 | 1.0 | 2.00 | 3.0 | 3.00 | 5.0 |
| 7 | 24.0 | 3.833333 | 0.868115 | 2.0 | 3.00 | 4.0 | 4.25 | 5.0 |
| 8 | 26.0 | 3.346154 | 1.324909 | 1.0 | 2.00 | 3.0 | 4.75 | 5.0 |
| 9 | 27.0 | 3.074074 | 1.639088 | 1.0 | 1.00 | 3.0 | 4.50 | 5.0 |
| 10 | 26.0 | 3.076923 | 1.440085 | 1.0 | 2.00 | 3.0 | 4.00 | 5.0 |
| 11 | 27.0 | 3.925926 | 1.141050 | 1.0 | 3.50 | 4.0 | 5.00 | 5.0 |
| 12 | 26.0 | 3.269231 | 1.250846 | 1.0 | 3.00 | 3.0 | 4.00 | 5.0 |
| 13 | 26.0 | 3.730769 | 1.041449 | 1.0 | 3.00 | 4.0 | 4.00 | 5.0 |
| 14 | 26.0 | 3.346154 | 1.383974 | 1.0 | 2.00 | 3.0 | 5.00 | 5.0 |
| 15 | 26.0 | 3.730769 | 1.401647 | 1.0 | 3.00 | 4.0 | 5.00 | 5.0 |
| 16 | 26.0 | 3.269231 | 1.401647 | 1.0 | 2.25 | 3.0 | 4.75 | 5.0 |
| 17 | 25.0 | 2.520000 | 1.084743 | 1.0 | 2.00 | 3.0 | 3.00 | 5.0 |
| 18 | 27.0 | 3.148148 | 1.166972 | 1.0 | 2.00 | 4.0 | 4.00 | 5.0 |
| 19 | 27.0 | 2.555556 | 1.187542 | 1.0 | 2.00 | 3.0 | 3.00 | 5.0 |
| 20 | 27.0 | 3.518519 | 1.155934 | 1.0 | 3.00 | 4.0 | 4.00 | 5.0 |
| 21 | 27.0 | 3.000000 | 1.441153 | 1.0 | 2.00 | 3.0 | 4.00 | 5.0 |
| 22 | 25.0 | 3.840000 | 1.027943 | 2.0 | 3.00 | 4.0 | 5.00 | 5.0 |
| 23 | 27.0 | 2.555556 | 1.187542 | 1.0 | 2.00 | 2.0 | 3.50 | 5.0 |
| 24 | 26.0 | 3.461538 | 1.475961 | 1.0 | 2.25 | 4.0 | 4.75 | 5.0 |
| 25 | 26.0 | 2.846154 | 1.461296 | 1.0 | 2.00 | 3.0 | 4.00 | 5.0 |
| 26 | 26.0 | 3.269231 | 1.250846 | 1.0 | 2.00 | 3.0 | 4.00 | 5.0 |
| 27 | 26.0 | 3.923077 | 0.934797 | 2.0 | 3.00 | 4.0 | 5.00 | 5.0 |
| 28 | 26.0 | 3.230769 | 1.242826 | 1.0 | 2.25 | 3.0 | 4.00 | 5.0 |
| 29 | 26.0 | 3.346154 | 1.294366 | 1.0 | 2.25 | 4.0 | 4.00 | 5.0 |
| 30 | 27.0 | 4.074074 | 0.916764 | 2.0 | 4.00 | 4.0 | 5.00 | 5.0 |

**Table 11: Estimated coefficients ($\hat{\beta}$) of the fixed effects of the mixed-effects model predicting explanation ratings.**

|  | Estimate | Std. Error | t value |
|---|---|---|---|
| (Intercept) | 3.31082 | 0.07517 | 44.044 |
| HighlightAVFALSE | -0.22321 | 0.08054 | -2.771 |
| TypeHuman | -0.15920 | 0.08038 | -1.981 |
| AVExperienceTRUE | 0.09922 | 0.05926 | 1.674 |
| HighlightAVFALSE :TypeHuman | 0.17437 | 0.11367 | 1.534 |

**Table 12: Variation of intercepts by scenario of the mixed-effects model predicting explanation ratings.**

|  | (Intercept) |
|---|---|
| Scenario 1 | 3.376390 |
| Scenario 2 | 3.259917 |
| Scenario 3 | 3.279770 |
| Scenario 4 | 3.327218 |

**Table 13: Estimated variation of the random effects, in this case just the Scenario, of the mixed-effects model predicting explanation ratings.**

| Groups | Name | Variance | Std.Dev. |
|---|---|---|---|
| Scenario | (Intercept) | 0.004611 | 0.0679 |
| Residual |  | 0.646100 | 0.8038 |

**Table 14: Analysis of Deviance Table (Type II Wald chi-square tests) with response variable of explanation ratings. Significance codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1**

|  | Chisq | Df | Pr(>Chisq) |
|---|---|---|---|
| HighlightAV | 5.6826 | 1 | * 0.01713 |
| Type | 1.6055 | 1 | 0.20513 |
| AVExperience | 2.8032 | 1 | . 0.09408 |
| HighlightAV:Type | 2.3529 | 1 | 0.12505 |

**Table 15: Estimated coefficients ($\hat{\beta}$) of the fixed effects of the mixed-effects model predicting trust ratings.**

|  | Estimate | Std. Error | t value |
|---|---|---|---|
| (Intercept) | 1.52881 | 0.50029 | 3.056 |
| MeanCEMARank | 0.35401 | 0.15279 | 2.317 |
| MeanHumanRank | 0.11929 | 0.15360 | 0.777 |
| PreExplanationFALSE | 0.11044 | 0.05158 | 2.141 |
| HighlightAVFALSE | 0.02302 | 0.05186 | 0.444 |
| AVExperienceTRUE | 0.33170 | 0.05248 | 6.321 |
| MeanCEMARank :MeanHumanRank | -0.04054 | 0.04477 | -0.906 |

**Table 16: Variation of intercepts by trust scales of the mixed-effects model predicting trust ratings.**

|  | (Intercept) |
|---|---|
|  | (Intercept) |
| Question 1 | 1.3028370 |
| Question 2 | 1.8321074 |
| Question 3 | 0.9450894 |
| Question 4 | 1.1313141 |
| Question 5 | 1.9791269 |
| Question 6 | 2.4642915 |
| Question 7 | 1.8762132 |
| Question 8 | 1.0921089 |
| Question 9 | 1.1362148 |

**Table 17: Estimated variation of the random effects, in this case just the trust scales, of the mixed-effects model predicting trust ratings.**

| Groups | Name | Variance | Std.Dev. |
|---|---|---|---|
| TrustQuestion | (Intercept) | 0.2787 | 0.5279 |
| Residual |  | 1.1299 | 1.0630 |

**Table 18: Analysis of Deviance Table (Type II Wald chi-square tests) with response variable of trust ratings. Significance codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1**

|  | Chisq | Df | Pr(>Chisq) |
|---|---|---|---|
| MeanCEMARank | 21.8790 | 1 | *** 2.904e-06 |
| MeanHumanRank | 0.0379 | 1 | 0.84563 |
| PreExplanation | 4.5852 | 1 | * 0.03225 |
| HighlightAV | 0.1970 | 1 | 0.65719 |
| AVExperience | 39.9520 | 1 | *** 2.603e-10 |
| MeanCEMARank :MeanHumanRank | 0.8202 | 1 | 0.36511 |

*D.3.2 Results of Statistical Analysis.* We present in the previous page the full results of our statistical analysis using the mixed-effects model for both the explanation ratings-related and trust ratings-related experiments In Tables 11 and 15, we report the estimated means and standard errors of the fixed effects. In Tables 12 and 16, we report the changes to the estimated mean of the intercepts due to the grouping by the random effects variables in the model. In Tables 13 and 17, we report the variation in the random effects variables. Finally, in Tables 14 and 18, we report the significance of each fixed effects variable on the predicted value of the outcome variable.

## D.4 ChatGPT Prompts and Replies

One possible additional baseline that has recently emerged is the family of large language models and their public-facing demonstrations such as ChatGPT [31]. While it is strongly debated whether these systems have a grasp on causality, nonetheless, their impressive fluency and apparent reasoning abilities should not be ignored.

Therefore, we have also explored the kind of explanations that ChatGPT can give by describing each scenario and then prompting it for an explanation. For all scenarios, tedious prompt editing and re-prompting were required to produce useful responses. We found that ChatGPT could identify causes correctly behind S1 and S4, and failed to explain S2 and S3. However, it is impossible to determine how and why specific causes are selected by ChatGPT. It is a closed system which changes frequently and its responses vary for the same prompts making it difficult to use as a baseline.

We give the prompts and the obtained responses from our exploration of ChatGPT's ability to explain behaviour in our scenarios. We used GPT-3.5 for these explanations.

*D.4.1 Scenario 1.* ChatGPT identifies a safety risk in the orange vehicle (non-ego) cutting in front of the blue vehicle (ego). After additional prompting, it gives further correct and plausible reasons to why the ego would change lanes left in this scenario.

> **Prompt:** There is a blue vehicle and an orange vehicle on a two-laned road. The orange vehicle is in the left lane 20 meters ahead of the blue vehicle which is in the right lane. The road is approaching a T-junction with an exit on the right. The orange vehicle changes lanes to the right and begins to slow down. In response, the blue vehicle changes lanes to the left. Why did the blue vehicle change lanes to the left?
>
> **Response:** The blue vehicle most likely changed lanes to the left to avoid a potential collision with the orange vehicle, which was slowing down and changing lanes to the right. Changing lanes to the left would give the blue vehicle more space and distance from the orange vehicle, reducing the risk of a collision as the road approaches the T-junction with an exit on the right. It's a defensive driving technique to maintain a safe distance from other vehicles and anticipate any potential hazards or obstacles on the road.
>
> **Prompt:** Besides avoiding a collision, what else could have prompted the blue vehicle to change lanes left?
>
> **Response:** There could be several other reasons why the blue vehicle changed lanes to the left besides avoiding a collision with the orange vehicle. Here are some possibilities:

> (1) Overtaking: The blue vehicle may have wanted to overtake another slower-moving vehicle in front of it, and changing lanes to the left provided a clear path to do so.
> (2) Obstruction: There may have been an obstruction in the right lane that the blue vehicle wanted to avoid, such as debris, a pothole, or a parked vehicle.
> (3) Turning left: The blue vehicle may have been planning to turn left at the T-junction, and changing lanes to the left put it in the correct position for the turn.
> (4) Personal preference: The driver of the blue vehicle may have simply preferred driving in the left lane for personal reasons, such as feeling more comfortable or having better visibility.
>
> However, without more context about the specific situation, it's impossible to know for sure why the blue vehicle changed lanes to the left. *[This last sentence is also by ChatGPT]*

*D.4.2 Scenario 2.* ChatGPT cannot identify the correct causes behind the blue vehicle's actions (ego) even when we give a large number of contextual clues. It gives completely incorrect guesses about why the blue vehicle might not give way to the orange vehicle (non-ego).

> **Prompt:** A blue vehicle is approaching a junction without traffic lights from the south on a secondary road. There is an orange vehicle coming from the left and a green vehicle coming from the right on a main road. The green vehicle is maintaining a high speed and is heading straight through the junction, and the orange vehicle is slowing down to a stop likely to turn left and waiting for the green vehicle to pass. In response, the blue vehicle turns right onto the main road without giving way to the orange vehicle. Why did the blue vehicle not give way?
>
> **Response:** Based on the scenario you described, the blue vehicle should have given way to the orange vehicle before turning right onto the main road. Failing to give way to the orange vehicle would be a violation of traffic laws in most jurisdictions, and could result in a collision or other dangerous situation.
> Without further information about the specific circumstances and the actions of the drivers involved, it's difficult to determine why the blue vehicle did not give way. However, some possible reasons could include the driver not being aware of the orange vehicle's presence or intention to turn left, or misjudging the distance and speed of the orange vehicle. It's important for drivers to always be aware of their surroundings and give way to other vehicles when required by traffic laws and common sense.

At this point, further prompting (e.g., telling ChatGPT explicitly that it is safe not to give way) does not change the responses of ChatGPT significantly.

*D.4.3 Scenario 3.* ChatGPT cannot identify the correct causes behind the blue vehicle's actions (ego) even when we give a large number of contextual clues. Further prompting (e.g., telling ChatGPT explicitly that it is safe not to give way) does not change the responses of ChatGPT significantly.

**Prompt:** There is a blue vehicle and an orange vehicle. The orange vehicle is in a two-laned roundabout and it has just changed lanes from the inner lane to the outer lane. The blue vehicle is on a road that enters the roundabout and it has to give way to the orange vehicle. On observing the orange vehicle change lanes, the blue vehicle decides to enter the roundabout without giving way to the orange vehicle. Why was it safe to enter the roundabout?

**Response:** As an AI language model, I do not have access to real-time information about the traffic conditions in your specific location. However, based on the scenario you have described, it is not safe for the blue vehicle to enter the roundabout without giving way to the orange vehicle, even if the orange vehicle had just changed lanes.

Failing to give way to a vehicle that has already entered a roundabout can cause a collision, particularly if the other vehicle is already in the driver's blind spot. It is important to remember that all drivers on the road have a responsibility to drive safely and follow traffic laws to prevent accidents. Therefore, it is recommended that the blue vehicle wait until the orange vehicle has passed

**Prompt:** However, the orange vehicle is also slowing down which indicates that it will exit the roundabout without crossing the path of the blue vehicle. Why was it safe to enter the roundabout for the blue vehicle now?

**Response:** If the orange vehicle is slowing down and is about to exit the roundabout without crossing the path of the blue vehicle, then it may be safe for the blue vehicle to enter the roundabout.

However, it is important to note that the blue vehicle should only enter the roundabout if it is safe to do so, and after ensuring that there are no other vehicles entering or exiting the roundabout that could potentially cross its path.

Additionally, the blue vehicle should always follow traffic laws and signals, and should be prepared to stop or yield to other vehicles if necessary.

Therefore, if the blue vehicle determines that it is safe to enter the roundabout and that it will not impede the flow of traffic or cause a potential collision, then it may proceed to enter the roundabout.

*D.4.4  Scenario 4.* Initially, ChatGPT assumes certain facts about the environment which are not true, but on some further clarifications, it can identify the correct causes behind the orange vehicle stopping at the junction.

**Prompt:** There is a waiting line of cars behind a red traffic light. Just behind this line of cars is a secondary road that enters the main road in a T-junction. There is a blue vehicle on the secondary road waiting to turn left onto the main road. An orange vehicle on the main road stops before the T-junction. In reaction, the blue vehicle turns left onto the main road and joins the waiting line of cars. Why did the blue vehicle turn left?

**Response:** Based on the information provided, it seems that the orange vehicle on the main road has stopped before the T-junction, potentially indicating that there is an obstruction or some other reason preventing it from continuing forward.

This may have prompted the blue vehicle on the secondary road to take advantage of the opportunity to turn left onto the main road and join the waiting line of cars, rather than waiting for the traffic light to turn green.

Alternatively, it is possible that the blue vehicle was already intending to turn left onto the main road, and the stopping of the orange vehicle simply provided a clear opening for the blue vehicle to make the turn safely. Without further information, it is difficult to determine the exact reason why the blue vehicle turned left onto the main road.

**Prompt:** There are no obstacles on the road so the orange vehicle could have just joined the line of waiting cars, thus blocking the way of the blue vehicle. Yet it did not. In light of this, why could the blue vehicle turn left?

**Respose:** If there were no obstacles on the road and the orange vehicle could have joined the line of waiting cars, but did not, it is possible that the orange vehicle driver intentionally stopped before the T-junction to let the blue vehicle turn left onto the main road. This could have been a courtesy gesture on the part of the orange vehicle driver, particularly if the blue vehicle had been waiting for a while or if the orange vehicle driver noticed that the blue vehicle was in a hurry.

Alternatively, it is possible that the blue vehicle had the right of way to turn left onto the main road and the orange vehicle driver recognized this, and therefore stopped to allow the blue vehicle to turn safely.

Without further information, it is difficult to determine the exact reason why the blue vehicle turned left onto the main road. However, it is possible that the actions of the orange vehicle driver played a role in allowing the blue vehicle to make the turn.