# A  Verbalisation Examples

In this section, we provide an example of how our implementation of AXIS for AD verbalises different elements of the environment context. We use scenario #2 to demonstrate how verbalisations look like in our experiments.

> The **road layout** verbalisation contains metadata and a textual summary of the roads, lanes, and their connections and priorities.

```
 1   - Metadata:
 2     - Traffic rules:
 3       - Vehicles drive on the right side of the road.
 4       - The maximum speed limit is 10 m/s.
 5     - Road layout:
 6       - The road layout consists of roads identified as Road(road
         ID).
 7       - Roads are made up of lanes identified as Road(road ID:
         lane ID).
 8       - Lanes are divided into left and right lanes.
 9       - Left lanes have positive lane IDs, while right lanes have
          negative lane IDs.
10     - Coordinate system:
11       - Distances are in meters.
12       - Coordinates are on a 2D Cartesian plane.
13       - Coordinates are written as [x, y].
14       - Angles are in radians in the range [-pi, pi].
15
16   - Road layout:
17     - Road(0): Road(0:1), Road(0:-1)
18     - Road(1): Road(1:1), Road(1:-1)
19     - Road(2): Road(2:1), Road(2:-1)
20     - Road(5): Road(5:1)
21     - Road(6): Road(6:-1)
22     - Road(7): Road(7:-1)
23     - Road(8): Road(8:-1)
24     - Road(9): Road(9:1)
25     - Road(10): Road(10:-1)
26     - Priorities:
27       - Road(6) has priority over: Road(8), Road(10), Road(9)
28       - Road(5) has priority over: Road(8)
29       - Road(9) has priority over: Road(8)
30
31   - Static objects:
32     - Building:
33       - Center: [-27. -16.]
34       - Boudnary: [[-49., -27.], [ -5., -27.], [ -5., -5.],
         [-49., -5.]]
```

> The **reward signal** verbalisation contains the reward components and their values.

```
 1   - Rewards:
 2       - Jolt: -0.03
 3       - Steering: -0.157
 4       - Curvature: -0.0
 5       - Time-to-Goal: 0.911
```

The **observation and option** verbalisation contains the times the positions and the velocities of each vehicle, subsampled at 3 Hz to decrease sequence lengths, as well as the option verbalisations which are written using the options' names and start and end times.

```
1   - Vehicle 0:
2     - Timesteps: [ 0, 3, 6, 9, 12, 15, 18, 21, 24, 27, 30, 33,
        36, 39, 42, 45, 48, 51, 54, 57, 60, 63, 66, 69, 72, 75, 78,
        81, 84, 87, 90, 93, 96, 99, 102, 105, 108, 111, 114, 117,
        121, 124, 127]
3     - Position: [[ 1.75, -53.79], [ 1.75, -52.6 ], [ 1.75,
        -51.36], [ 1.75, -50.07], [ 1.75, -48.74], [ 1.75, -47.39],
        [ 1.75, -46. ], [ 1.75, -44.59], [ 1.75, -43.17], [ 1.75,
        -41.72], [ 1.75, -40.27], [ 1.75, -38.8 ], [ 1.75, -37.32],
        [ 1.75, -35.83], [ 1.75, -34.33], [ 1.75, -32.83], [ 1.75,
        -31.32], [ 1.75, -29.81], [ 1.75, -28.29], [ 1.75,
        -26.77], [ 1.75, -25.25], [ 1.75, -23.73], [ 1.75, -22.2 ],
        [ 1.75, -20.67], [ 1.75, -19.15], [ 1.75, -17.69], [ 1.75,
        -16.3 ], [ 1.75, -14.95], [ 1.75, -13.66], [ 1.75,
        -12.44], [ 1.75, -11.32], [ 1.75, -10.31], [ 1.75, -9.42],
        [ 1.75, -8.66], [ 1.75, -8.01], [ 1.75, -7.47], [ 1.75,
        -7.06], [ 1.75, -6.76], [ 1.75, -6.59], [ 1.75, -6.53], [
        1.75, -6.53], [ 1.75, -6.53], [ 1.75, -6.53]]
4     - Speed: [7.5 , 7.87, 8.19, 8.46, 8.69, 8.89, 9.05, 9.19,
        9.31, 9.41, 9.5 , 9.58, 9.64, 9.7 , 9.74, 9.78, 9.81, 9.84,
        9.86, 9.88, 9.9 , 9.92, 9.93, 9.94, 9.77, 9.37, 8.9 ,
        8.65, 8.25, 7.75, 7.04, 6.27, 5.5 , 4.74, 3.97, 3.2 , 2.43,
        1.66, 0.89, 0.12, 0. , 0. , 0. ]
5     - Steering: KeepStraight[0-129]
6     - Acceleration: Accelerate[0-73], Decelerate[73-121],
        MaintainSpeed[121-129], Accelerate[129-129]
7     - Macro actions (as macro[from-to]): [FollowLane[0-72],
        GiveWay[73-120], Stop[121-128], GiveWay[129-129]]
8     - Lane sequence (as Road(...)[from-to]): [Road(2:-1)[0-127]]
9
10  - Vehicle 1:
11    - Timesteps: [ 0, 3, 6, 9, 12, 15, 18, 21, 24, 27, 30, 33,
        36, 39, 42, 45, 48, 51, 54, 57, 60, 63, 66, 69, 72, 75, 78,
        81, 84, 87, 90, 93, 96, 99, 102, 105, 108, 111, 114, 117,
        120, 123, 126]
12    - Position: [[16.34, 1.75], [15.36, 1.75], [14.41, 1.75],
        [13.51, 1.75], [12.65, 1.75], [11.84, 1.75], [11.09, 1.75],
        [10.39, 1.75], [ 9.74, 1.75], [ 9.15, 1.75], [ 8.66,
        1.75], [ 8.3 , 1.75], [ 8.03, 1.75], [ 7.84, 1.75], [ 7.7 ,
        1.75], [ 7.6 , 1.75], [ 7.53, 1.75], [ 7.48, 1.75], [
        7.45, 1.75], [ 7.42, 1.75], [ 7.4 , 1.75], [ 7.39, 1.75], [
        7.38, 1.75], [ 7.37, 1.75], [ 7.37, 1.75], [ 7.37, 1.75],
        [ 7.36, 1.75], [ 7.36, 1.75], [ 7.36, 1.75], [ 7.36, 1.75],
        [ 7.36, 1.75], [ 7.36, 1.75], [ 7.36, 1.75], [ 7.36,
        1.75], [ 7.36, 1.75], [ 7.36, 1.75], [ 7.36, 1.75], [ 7.36,
        1.75], [ 7.36, 1.75], [ 7.35, 1.75], [ 7.25, 1.75], [
        7.06, 1.75], [ 6.79, 1.75]]
13    - Speed: [6.61, 6.43, 6.18, 5.88, 5.56, 5.21, 4.85, 4.52, 4.2
        , 3.69, 2.93, 2.18, 1.57, 1.12, 0.81, 0.58, 0.41, 0.3 ,
        0.21, 0.15, 0.11, 0.08, 0.06, 0.04, 0.03, 0.02, 0.01, 0.01,
        0.01, 0.01, 0. , 0. , 0. , 0. , 0. , 0. , 0. , 0.01, 0.01,
        0.15, 0.68, 1.12, 1.48]
14    - Steering: KeepStraight[0-129]
15    - Acceleration: Decelerate[0-62], MaintainSpeed[62-117],
        Accelerate[117-129]
16    - Macro actions (as macro[from-to]): [GiveWay[0-61], Stop
        [62-116], GiveWay[117-129]]
17    - Lane sequence (as Road(...)[from-to]): [Road(1:1)[0-126]]
```

# B  Algorithm Prompt Templates

This section gives the verbatim text of the prompts used in AXIS. There are five prompt templates: `system`, `context`, `interrogation`, `explanation`, and `final`. The prompt templates are also available in their original form in the codebase, under the *'data/igp2/prompts'* directory.

**Placeholders in system prompt templates:**  Note, when no simulation is used, the <N_MAX> variables is not given, and the prompt is slightly changed to remove reference to simulation rounds.

- **<N_MAX>**: The maximum number of interrogation-synthesis rounds $N_{\max}$ in AXIS.
- **<EXPLANATION_STYLE>**: The instructions on the stylistic constraints of the explanation. In out experiments this was set to: *"Do not include raw state or action arrays. Do not include explicit references to road, lane, and intersection IDs. Refer to casual relationships."*
- **<COMPLEXITY>**: The linguistic complexity of the explanation. We use two options:
  1. Simple: Focusing on brevity and simplicity; *"Your response must be as short and concise as possible. You can only include the absolute most important information."*
  2. Complex: Focusing on clarity and concision; *"Your response must be concise and clear. Include only important information."*

---

The system prompt template of the LLM used in the AXIS agent.

```
1  You write helpful explanations based on a multi-round sequence
       of dialog over {n_max} rounds. Each round consists of two
       Stages:
2
3  1. Interrogation Stage: You interrogate a simulator via a query
        to obtain additional counterfactual information about the
       world;
4  2. Explanation Stage: You synthesize all information available
       to you into a concise explanation.
5
6  {explanation_style}
7  {complexity}
8  Do not give an explanation until explicitly prompted.
9  End all of your responses with '<|endoftext|>'.
```

---

**Placeholders in context prompt templates:**  Note, some of these placeholders may not appear in versions of the prompt that are used for baselines, i.e. not simulation is used, or when no initial context is given, i.e. when the LLM needs to build an understanding of the scenario on its own.

- **<CONTEXT>**: The verbalised context elements as a single contiguous string (e.g. static road layout $V_{env}$, observations $V_o$, etc.).
- **<OCCLUSIONS>**: Optional text, given when an occlusion may be present in the scenario; *"Note, the context observations and actions may be incomplete as the scenario could potentially contain occluded vehicles."*
- **N_MAX**: The maximum number of interrogation-synthesis rounds $N_{\max}$ in AXIS.
- **<QUERY_DESCRIPTION>**: The description of the interrogation queries as a string.
- **<QUERY_TYPE_DESCRIPTION>**: The data type descriptions of each argument of the interrogation queries as a string (e.g. time $\in \mathbb{N}$, actions $\in \mathcal{M}^*$, etc.).
- **<USER_PROMPT**: The user prompt $p_u$.

> Initial context prompt template used to setup the AXIS explanation generation process.

```
1  You explain why a vehicle took an action in a particular
       scenario. You learn about causes and effects by asking me
       queries over {n_max} rounds, proposing changes to what was
       observed in the scenario and observing what happened
       differently thereby collecting information. {occlusions}
       You are given the following initial context about the
       observations and actions of the vehicles in the scenario:
2  {context}
3
4  You must choose one from the list of valid queries. Each query
       has a name and a list of keyword arguments following the
       syntax: 'query(arg1=value1, arg2=value2, ...)'. The list of
        valid queries and their meanings are:
5  {query_descriptions}
6
7  The keyword arguments must have the following types:
8  {query_type_descriptions}
9
10  You need to collect information to explain the following
       question:
11  {user_prompt}
```

**Placeholders in interrogation prompt templates:**

- **<N>**: The current round number $N$ of the interrogation-synthesis process.

> The interrogation prompt template used to elicit an interrogation query from the LLM.

```
1  This is round {n}.
2  This is the Interrogation stage.
3
4  Choose a query from the list of valid queries. Specify concrete
       values for the keyword arguments. Reply only with a single
       query. Do not add anything else in your answer.
5  You may also answer DONE if you do not need more information.
```

**Placeholders in explanation prompt templates:**

- **<N>**: The current round number $N$ of the interrogation-synthesis process.
- **<CONTEXT>**: The verbalised context elements of the simulation results as a single contiguous string. See context prompt template above for further details.
- **<OCCLUSIONS>**: Optional text, given when an occlusion may be present in the scenario. See context prompt template above for details.
- **<EXPLANATION_STYLE>**: The instructions on the stylistic constraints of the explanation. See system prompt template above for details.
- **<COMPLEXITY>**: The linguistic complexity of the explanation. See system prompt template above for details.
- **<USER_PROMPT**: The user prompt $p_u$.

> The explanation prompt template used to elicit an intermediate explanation from the LLM.

```
1  This is round {n}.
2  This is the Explanation Stage.
3
4  After running the simulation for your query, the following was
       observed:
5  {context}
6
7  I want you to give an explanation to the below question. {
       occlusions}
8  {explanation_style}
9  {complexity}
10 {user_prompt}
```

**Placeholders in final prompt templates:**

- **<USER_PROMPT**: The user prompt $p_u$.
- **<OCCLUSIONS>**: Optional text, given when an occlusion may be present in the scenario. See context prompt template above for details.
- **<EXPLANATION_STYLE>**: The instructions on the stylistic constraints of the explanation. See system prompt template above for details.
- **<COMPLEXITY>**: The linguistic complexity of the explanation. See system prompt template above for details.

> The final prompt template used to synthesis the final explanation.

```
1  Generate your final explanation using all previous information
       in response to the question: {user_prompt}
2  {occlusions}
3  {explanation_style}
4  {complexity}
```

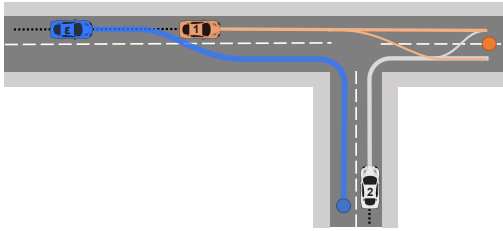## C   Scenarios and User Prompts for AD Evaluation

The 'longtable' below provides detailed descriptions of all ten scenarios as well as the user prompts which were used for experiments with AXIS. Illustrations of all scenarios (except the irrational ones) are given in Figures 1 and 2. Note, in the codebase, the scenarios are numbered from zero, but we number them in the paper from one. Also note that these descriptions were not the ones used for the perceived correctness metric, which are instead given in Appendix D.

For each scenario definition, we describe the road layout, the scenario itself, and the scenario category. The layout is given as a textual description of the road, including the number of lanes, their direction, and the junctions. The scenario is a textual description of the vehicles, their actions, and the environment. The category is a label that indicates the type of scenario, and can be rational, irrational, or occlusion.
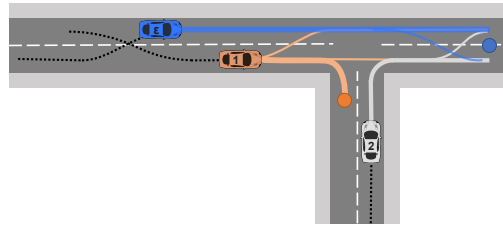
| ID | Description and Category | User Prompts |
|----|--------------------------|--------------|
| 1 | **Layout:** Road with two parallel lanes and a T-junction coming up ahead. **Scenario:** Two vehicles. The queried vehicle starts in the left lane then changes lanes right to prepare for an exit right. Vehicle 1 in the left lane continues, maintaining its speed. **Category:** Rational | **Prompt:** *Why did vehicle 0 change lanes right?* (Time: 80) **Prompt:** *Why didn't vehicle 0 go straight in the junction?* (Time: 180) |

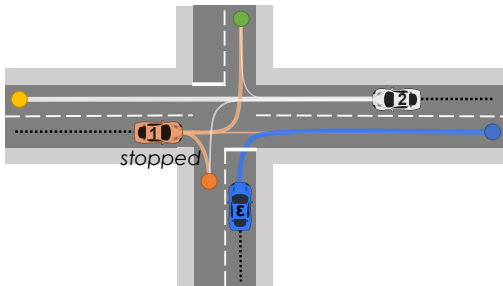| 2 | **Layout:** Same as scenario #1. <br> **Scenario:** Three vehicles. The queried vehicle starts in the right lane, when vehicle 1 cuts in front of it from the left lane and begins to slow down rapidly, indicating its intention to exit right at the junction. In response, the queried vehicle changes lanes left to continue to its goal unimpeded and more safely. Vehicle 2 is giving way at the T-junction to both the queried vehicle and vehicle 1. <br> **Category:** Rational | **Prompt:** *Why did vehicle 0 not go straight?* (Time: 80) <br> **Prompt:** *Why did vehicle 0 change lanes left?* (Time: 80) <br> **Prompt:** *What if vehicle 1 hadn't changed lanes right?* (Time: 80) |
|---|---|---|
| 3 | **Layout:** Four-way crossroads between a main road and two lower-priority roads. Each road has two lanes, with one lane in each driving direction. The roads are perpendicular to one another. <br> **Scenario:** Three vehicles. The queried vehicle is approaching the crossroads from a low-priority road. Vehicle 1 is on the main road, slowing down to a stop as it approaches the crossroads. Vehicle 2 is going down the main road in the opposite direction of vehicle 1 while maintaining its speed close to the speed limit. The queried vehicle realises that vehicle 1 is trying to turn left and is giving way to vehicle 2 on the main road. The queried vehicle can use this opportunity to turn right early without waiting for vehicle 1 to turn. <br> **Category:** Rational | **Prompt:** *What would have happened if vehicle 1 had gone through the junction instead of stopping?* (Time: 160) <br> **Prompt:** *Why did vehicle 0 not stop?* (Time: 160) <br> **Prompt:** *What would have happened if vehicle 0 had stopped?* (Time: 160) |
| 4 | **Layout:** Roundabout with three exits and two lanes. <br> **Scenario:** Two vehicles. The queried vehicle is approaching the roundabout from the south. Vehicle 1 is in the roundabout in the inner lane at the start but changes lanes to the left. The queried vehicle in reaction enters the roundabout without giving way to vehicle 1 as it infers the intention of vehicle 1 to exit the roundabout which was indicated by vehicle 1's left lane change. <br> **Category:** Rational | **Prompt:** *Why did vehicle 0 not stop at the roundabout?* (Time: 140) <br> **Prompt:** *Why did vehicle 0 not stop to give way at the roundabout?* (Time: 140) |
| 5 | **Layout:** T-junction followed by a four-way crossroads. <br> **Scenario:** Five vehicles. Two vehicles are waiting in line at the four-way crossroads at a traffic light. Vehicle 3 is approaching behind them. Vehicle 4 is passing through the four-way crossroads in the opposite direction of the waiting cars. The queried vehicle is approaching from the T-junction and is aiming to merge behind the waiting line of cars. Vehicle 3 sees this and slows down to stop, leaving a gap for the queried vehicle to merge. The queried vehicle realises this and uses the gap to merge behind the waiting line of cars. <br> **Category:** Rational | **Prompt:** *What would vehicle 0 do if vehicle 3 had gone through the junction instead of stopping?* (Time: 140) <br> **Prompt:** *Why did vehicle 0 not stop at the junction?* (Time: 140) <br> **Prompt:** *Why did vehicle 0 merge instead of giving way?* (Time: 140) |

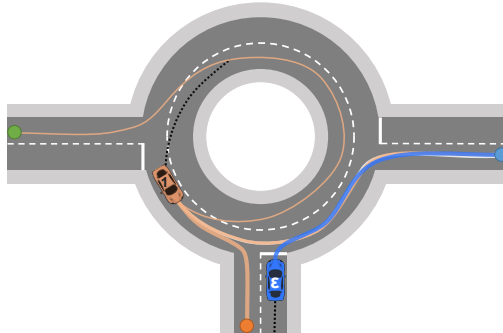| | | |
|---|---|---|
| 6 | **Layout:** Same as scenario #3. **Scenario:** The scenario is the same as scenario #3, however vehicle 1 after slowing down decided to speed back and head straight which results in a collision with the turning queried vehicle, which thought vehicle 1 is about to turn left. **Category:** Irrational | **Prompt:** *Why did vehicle 0 collide?* (Time: 160) **Prompt:** *How could have vehicle 0 avoided the collision?* (Time: 160) |
| 7 | **Layout:** Same as scenario #2. **Scenario:** The scenario is the same as scenario #2, however vehicle 1 is changing lanes back and forth on the main road. In response, the queried vehicle stays in the right lane instead of changing lanes to the left and keeps some distance from vehicle 1. **Category:** Irrational | **Prompt:** *Why did vehicle 0 not change lanes to the left?* (Time: 100) |
| 8 | **Layout:** Same as scenario #2. **Scenario:** Three vehicles. There is a parked vehicle 2 in the left lane of the two-lane main road. Behind it is vehicle 1 and then the queried vehicle. Vehicle 1 is blocking the view of vehicle 2 from the perspective of the queried vehicle. However, the queried vehicle observes vehicle 1 changing lanes to the right which would otherwise not be rational unless vehicle 2 were present, so to avoid the inferred parked vehicle, the queried vehicle also changes lanes. **Category:** Occlusion | **Prompt:** *Why did vehicle 0 not go straight instead of changing lanes to the right?* (Time: 100) **Prompt:** *Why did vehicle 0 change lanes right?* (Time: 100) |
| 9 | **Layout:** Same as scenario #3 with the addition of a large building that blocks the view of the queried vehicle to the left. **Scenario:** Three vehicles. The queried vehicle is approaching from a low-priority road and can observe vehicle 1 on the main road from the right, coming to a rolling stop. However, the queried vehicle cannot see what is on the left because the building blocks its view. Still, from the actions of vehicle 1, it is inferred that there is an oncoming vehicle on the main road from the left. Therefore, the queried vehicle stops to give way to the inferred vehicle. **Category:** Occlusion | **Prompt:** *Why did vehicle 0 stop to give way?* (Time: 130) **Prompt:** *Why would vehicle 0 not turn when there is no vehicle to give way to?* (Time: 130) |
| 10 | **Layout:** Same as scenario #3 with the addition of a large building that blocks the view of the queried vehicle to the right. **Scenario:** Three vehicles. The queried vehicle is approaching the main road from a low-priority road. It observes vehicle 1 on the main road from the left coming to a full stop for a longer amount of time. As vehicle 1 does not seem to want to turn and there is no vehicle appearing from the right side of the main road, the queried vehicle infers that the only rational reason for vehicle 1 to stop would be that there is a vehicle blocking its path occluded by the building on the right. Once the queried vehicle infers this, it decides to turn left, using the gap that vehicle 1 has left at the crossroads. **Category:** Occlusion | **Prompt:** *Why did vehicle 1 stop?* (Time: 140) **Prompt:** *Why did vehicle 0 not stop to give way?* (Time: 140) |

(**#1**) The blue queried vehicle changes lanes left, slows down, and turns right, because its goal is to exit right at the junction. Other vehicles have a minor safety-related teleological effect, but otherwise do not have a mechanistic causal effect on the queried vehicle.

(**#2**) The orange vehicle in front of the blue queried vehicle 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 queried vehicle decides to change lanes as it is heading straight.
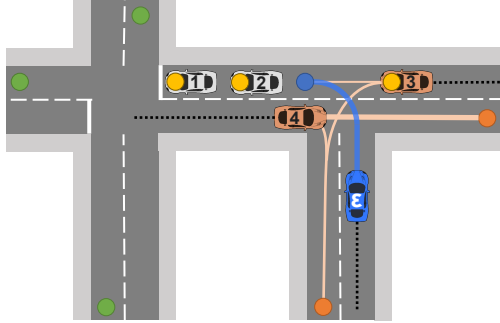
(**#3**) The queried vehicle is turning right but must give way. It observes vehicle 1 on the left stopping. This is only rational if it is trying to turn left and is giving way for the oncoming gray vehicle. The queried vehicle can use this to enter the road earlier.
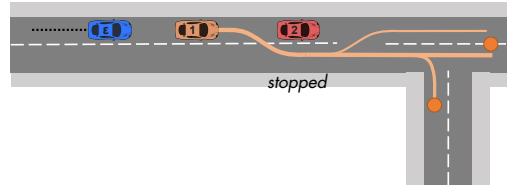
(**#4**) The queried vehicle observes the orange vehicle changing lanes to the right. This is only rational if the orange vehicle is leaving the roundabout at the next exit. The queried vehicle can therefore enter the roundabout faster without waiting to give way.
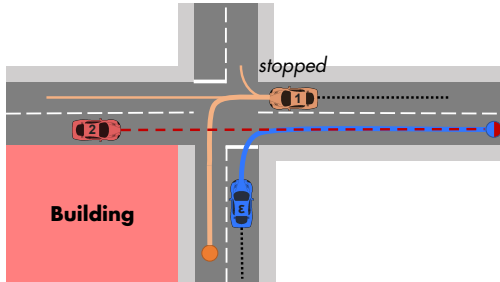
Figure 1: Solid lines show predicted trajectories of other vehicles with their thickness corresponding to predicted probability. Black dotted lines denote observations.
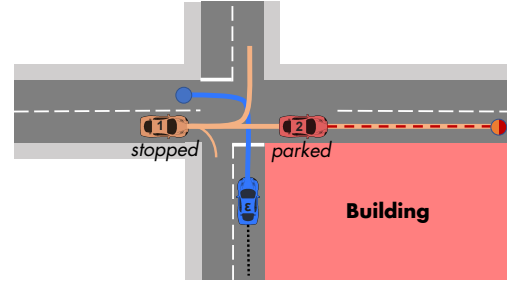
(**#5**) Vehicle 3 is slowing down to stop. Once vehicle 4 drives past as indicated by its maintained high speed, the stopping of vehicle 3 stays rational only if it is to allow the queried vehicle to merge without waiting for vehicle 4 to pass.

(**#8**) Vehicle 1 is slowing down then changes lanes while blocking the queried vehicle's view on vehicle 2. If vehicle 1's goal is to go straight, then its actions are irrational, so the queried vehicle infers that there is either an occluded vehicle or vehicle 1 must be turning right. To avoid any collisions with vehicle 2, the queried vehicle changes lanes right.

(**#9**) Vehicle 1 has priority but stops on the main road, from which the queried vehicle can infer that there is likely an occluded vehicle behind the building, and vehicle 1 is giving way to it, otherwise its actions are irrational. To avoid potential collisions, the queried vehicle decides to stop.

(**#10**) Vehicle 1 has priority on the main road but stops and waits at the intersection. After a while, queried vehicle infers that vehicle 1's actions are only rational if there is a stopped occluded vehicle behind the building, so queried vehicle decides to turn left using the gap left by vehicle 1.

Figure 2: Circles indicate agents' goals. Solid lines represent the predicted trajectories of other vehicles, with their thickness indicating the predicted probability. Black dotted lines denote observations, while red dashed lines denote the actual trajectory of the occluded agent.

## D   LLM-as-a-Judge Details

This section presents the hyperparameters used for sampling for LLMs, as well as the verbatim text of the prompts used in evaluating the explanations. There are three evaluation prompt templates: `actionability`, `correctness`, and `subjective preferences`. The prompt templates are also available in their original form in the codebase, under *'data/igp2/evaluation/'*.

We used vLLM[1] for local inference and vLLM's OpenAI-style API to access all other proprietary methods. The hyperparameters used for sampling are those provided by vLLM, most importantly, we used a temperature of 1.0, a top-$p$ sampling parameter of 1, and we set the seed to 28.

**Placeholders in actionability evaluation prompt templates**

- **<SCENARIO>**: The scenario context verbalisation (e.g. static road layout $V_{env}$, observations $V_o$, etc.).
- **<EXPLANATION>**: The AXIS-generated final explanation $e$. This may be omitted for the *NoExplanation* baseline.
- **<GOALS>**: A list of four possible goals the agent may be pursuing, printed in a randomised order to prevent sequence effects.
- **<MANEUVERS>**: Same as above, for options.

---

Actionability evaluation prompt template.

```
1   You are given the Scenario and the Explanation below.
2
3   Scenario:
4   {scenario}
5
6   Explanation:
7   {explanation}
8
9   Predict the final Goal of Vehicle 0 and what its next Maneuver
        will be.
10  For the Goal, you must choose one option from the following
        list 4 possible Goals:
11  {goals}
12  For the Maneuver, you must choose one option from the following
         list of 4 possible Maneuvers:
13  {maneuvers}
14
15  For both Goal and Maneuver, your response must be a single
        number indicating which Goal and Maneuver you have chosen.
16  Give your final answer in the format below.
17
18  $ Goal: <your chosen Goal>
19  $ Maneuver: <your chosen Maneuver>
```

---

**Placeholders in correctness evaluation prompt templates**

- **<SCENARIO>**: The scenario context verbalisation. See the actionability template for details.
- **<DESCRIPTION>**: The expert description of the scenario as baseline for the LLM to compare the explanation against. For each scenario ID the descriptions are:
  1. Vehicle 0 changed lanes left because it had to turn right at the junction. Vehicle 1 had no impact on its decision-making.
  2. Vehicle 0 changed lanes left when Vehicle 1 cut in front of it from the left lane and began to slow down. Vehicle 0 inferred that Vehicle 1 was trying to turn right at the junction, which would have slowed down Vehicle 0's path, so Vehicle 0 overtook Vehicle 1.

---

[1] https://docs.vllm.ai/en/latest/

3. Vehicle 0 did not stop because it saw Vehicle 1 coming to a halt at the intersection. Vehicle 0 inferred that Vehicle 1 was going to give way to Vehicle 2, which was maintaining its velocity. Vehicle 0 saw this as an opportunity to turn right without stopping to give way to Vehicle 1, which was more efficient and still safe.

4. Vehicle 0 entered the roundabout without giving way to Vehicle 1, which was in the roundabout and thus had priority, because Vehicle 1 changed to the outer lane in the roundabout. From this, Vehicle 0 could infer that Vehicle 1 was going to exit the roundabout. Vehicle 0 therefore could enter the roundabout because it was safe.

5. Vehicle 0 turned left early at the junction because it saved Vehicle 3 from coming to a stop and leaving a gap for Vehicle 0 to merge onto the road. If Vehicle 3 had not left the gap and had continued to drive straight instead, then Vehicle 0 would have had to stop and wait to merge until the lane was clear.

6. Vehicle 0 collided with Vehicle 1 because Vehicle 1 was acting irrationally. Vehicle 1 first came to a stop at the junction, indicating that it was giving way to Vehicle 2 while waiting to turn left, but then it decided to go straight instead. Vehicle 0 was already in the process of turning right and could not stop in time to avoid the collision. Vehicle 0's decision to turn right was based on the assumption that Vehicle 1 would give way to it, which was not the case.

7. Vehicle 0 did not change lanes because it was the most efficient way to reach its goal. In addition, the behaviour of Vehicle 1, which was erratic and unpredictable, made it necessary to maintain a safe distance from Vehicle 1. Vehicle 2's actions did not influence Vehicle 0 as Vehicle 2 was giving way at a junction.

8. Vehicle 0 changed lanes right when it saw Vehicle 1 change lanes right. Normally, Vehicle 1's actions would have been irrational, but Vehicle 0 inferred from the actions of Vehicle 1, that Vehicle 2 must be blocking the way even though Vehicle 2 was not observed initially.

9. Vehicle 0 stopped to turn right at the junction because it saw Vehicle 1 coming to a stop at the junction. Vehicle 0 inferred that Vehicle 1 was going to give way to an occluded Vehicle 2. Vehicle 0 safely waited until Vehicle 2 passed before turning right.

10. Vehicle 0 turned left without waiting for Vehicle 1 to pass, even though Vehicle 1 had priority, because Vehicle 1 came to a stop at the junction. Vehicle 0 inferred that the path of Vehicle 1 was blocked by an occluded Vehicle 2. Vehicle 0, therefore, could turn left safely without stopping, which was more efficient.

- **<QUESTION>**: The user prompt $p_u$.
- **<EXPLANATION>**: The AXIS-generated final explanation $e$.

---

Perceived correctness evaluation prompt template.

```
1  You are given a driving Scenario with a corresponding true
      Description of what happens in the Scenario. You are also
      given an Explanation to a Question about the actions and
      behaviour of a vehicle in the Scenario.

2

3  Scenario:
4  {scenario}

5

6  Description:
7  {ground_truth}

8

9  Question:
10 {question}

11

12 Explanation:
13 {explanation}

14

15 Determine whether the Explanation is correct.
16 Your response must be an integer score on a scale of 1 to 5,
      where 1 is completely incorrect, 2 is mostly incorrect i, 3
       is partially correct, 4 is mostly correct, and 5 is
      completely correct.
```

```
17   You response must also include a short and concise Reasoning of
         your chosen Score.
18   You must give your final answer in the format below.
19
20   $ Score: <your Correctness score>
21   $ Reasoning: <your Reasoning>
```

**Placeholders in subjective preferences evaluation prompt templates**

- **<SCENARIO>**: The scenario context verbalisation. See the actionability template for details.
- **<QUESTION>**: The user prompt $p_u$.
- **<EXPLANATION>**: The AXIS-generated final explanation $e$.

Subjective preferences evaluation prompt template.

```
1    You are given an Explanation to a Question about the actions
         and behaviour of a vehicle in a given context Scenario.
2
3    Scenario:
4    {scenario}
5
6    Question:
7    {question}
8
9    Explanation:
10   {explanation}
11
12   Evaluate how well the Explanation answers the Question given
         the Scenario by answering how much you agree or disagree
         with the following statements:
13   - SufficientDetail: This explanation of why the self-driving
         car behaved as it did has sufficient detail.
14   - Satisfying: This explanation of why the self-driving car
         behaved as it did is satisfying.
15   - Complete: This explanation of why the self-driving car
         behaved as it did seems complete.
16   - Trust: This explanation lets me judge when I should trust and
          not trust the self-driving car.
17
18   For each statement, your response must be an integer score on a
         scale of 1 to 5, where 1 is strongly disagree, 2 is
         somewhat disagree, 3 is neither agree nor disagree, 4 is
         somewhat agree, and 5 is strongly agree.
19   Give your answer in the following format.
20
21   $ SufficientDetail: <your score for SufficientDetail>
22   $ Satisfying: <your score for Satisfying>
23   $ Complete: <your score for Complete>
24   $ Trust: <your score for Trust>
```
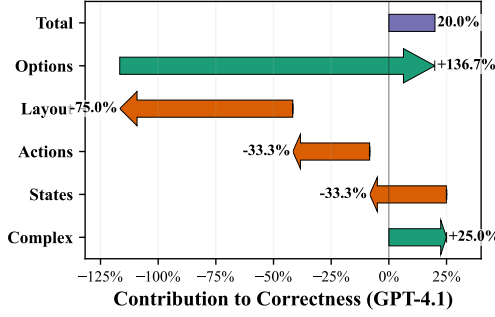
# E  Shapley Value Analysis

This section describes the methodology and results of the Shapley value analysis presented in the main paper. As described in the main paper in Section 4, we use 5 features for which we calculate their Shapley value contributions to the final correctness score of the explanation. The features are verbalised aspects of the environment context, and they are:

- **Road Layout**: The semantic road layout, driving rules, and metadata (e.g. units) verbalised as $V_{env}$.
- **States**: The state vector sequences of the verbalised observations $V_o$, i.e. position, velocity, and heading.
- **Actions**: The action vector sequences of verbalised observations $V_o$, i.e. acceleration and steering values.
- **Options**: The high-level options, e.g. turn right, give way, $V_m$.
- **Linguistic Complexity**: The linguistic complexity of the explanation, i.e. concise or complex.
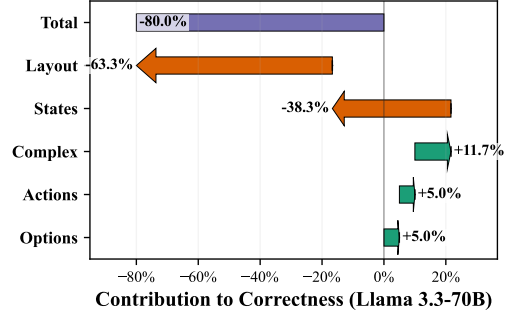
Formally, let us denote the set of features above as $V$. The Shapley value $\psi_i$ of a feature $V_i$ is defined as the average marginal contribution of that feature to all possible subsets of features, as given by:

$$\psi_i = \sum_{S \subseteq V \setminus \{V_i\}} \frac{|S|!(|V| - |S| - 1)!}{|V|!} \left( v(S \cup \{V_i\}) - v(S) \right), \tag{1}$$

where $v(S)$ is the value — in our case the perceived correctness — of the subset of features $S$, and $|S|$ is the number of features in the subset. The Shapley value is a measure of the contribution of each feature to the overall perceived correctness of the explanation. We calculate $\psi_i$ for each feature $V_i$ by evaluating the correctness of explanations generated with all possible subsets of features, then using Equation (1) to calculate the Shapley value for each feature.

(a) Scenario #3 (rational)- GPT-4.1

(b) Scenario #3 (rational) - Llama 3.3-70B

(c) Scenario #7 (irrational) - GPT-4.1

(d) Scenario #7 (irrational) - Llama 3.3-70B

(e) Scenario #8 (occlusion) - GPT-4.1

(f) Scenario #8 (occlusion) - Llama 70B

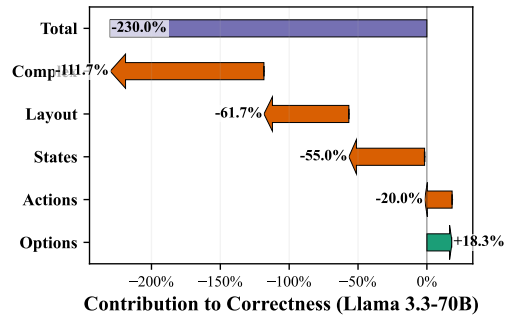Figure 3: Shapley value analysis comparing the contribution of different features to explanation performance as measured by correctness across three scenarios, #3, #7, and #8. Left plots show GPT-4.1 results, while right plots show Llama 3.3-70B results.

# F  Further Qualitative Examples

Results for scenario #1 — Question: *Why did vehicle 0 changes lane right?*

- **GPT-4.1:** Vehicle 0 changed lanes to the right because this maneuver was necessary to reach its goal. When instructed to avoid the lane change and simply continue straight, vehicle 0 was unable to complete its route, demonstrating that the lane change was essential for making an upcoming turn or aligning with the required path. This decision was not caused by the presence of vehicle 1 but was instead dictated by the route requirements and vehicle 0's intended destination.

- **DeepSeek-V3:** Vehicle 0 changed lanes right to maintain efficient travel speed and avoid being blocked by Vehicle 1. When both vehicles were in the same lane, Vehicle 0 had to slow down significantly, hindering its progress. The lane change allowed Vehicle 0 to continue unimpeded, demonstrating that its primary motivation was to optimize forward movement rather than respond to Vehicle 1's specific position. This strategic maneuver prevented unnecessary deceleration while maintaining safe operation.

- **Llama 3.3-70B:** Vehicle 0 changed lanes right to reach its goal, and this action was not directly influenced by vehicle 1's actions, suggesting it was a deliberate maneuver to achieve its goal.

- **o4-mini:** Vehicle 0's shift to the right lane was purely preparatory for its upcoming right-turn. By moving over well before the junction, it secured the correct position to execute the turn; without that early lane change it would miss the turn and fail to reach its goal. This maneuver occurred even when no other vehicles were present, confirming it wasn't a reaction to traffic.

- **DeepSeek-R1:** Vehicle 0 changed lanes right to preemptively resolve a conflict between its route goals and Vehicle 1's trajectory. The maneuver served two causal purposes: 1. **Obstacle Avoidance**: Vehicle 1's steady-speed parallel path would have forced Vehicle 0 to decelerate (observed in counterfactuals), reducing efficiency. 2. **Route Optimization**: The lane change aligned with Vehicle 0's inherent path objectives, avoiding later sharp turns/speed penalties required when no obstacle existed (removal scenario). By acting early, Vehicle 0 balanced immediate conflict mitigation with long-term trajectory efficiency.

Results for scenario #3 — Question: *Why did vehicle 0 not stop?*

- **GPT-4.1:** Vehicle 0 did not stop because it was able to coordinate its movement with the other vehicles present. It yielded as necessary but was not forced to a complete stop, as vehicle 1's and vehicle 2's trajectories allowed for safe passage through the conflict area. Removing either vehicle or forcing them to stop resulted in vehicle 0 becoming stuck or delayed, showing that its ability to keep moving depended on their continuous movement and appropriate yielding. The timing and presence of the other vehicles enabled vehicle 0 to proceed without stopping.

- **DeepSeek-V3:** Vehicle 0 did not stop because: 1. **Physical Constraints**: It was too close to the lane end to safely stop without violating traffic rules. 2. **Traffic Flow Priority**: Maintaining movement was more efficient than a full stop, given Vehicle 1's consistent speed and trajectory. 3. **Safe Interaction**: The GiveWay action allowed sufficient slowing to yield to Vehicle 1 while avoiding collision. 4. **Cooperative Behavior**: Both vehicles adapted smoothly—Vehicle 0 yielded minimally while Vehicle 1 maintained speed, optimizing overall junction efficiency. Stopping would have been unnecessarily disruptive without improving safety.

- **Llama 3.3-70B:** Vehicle 0 didn't stop as it had a turning intention and was yielding to other vehicles, but when it was clear to proceed, it continued to turn, prioritizing its goal over stopping.

- **o4-mini:** Vehicle 0 slowed to yield to conflicting flows but never hit zero because each opposing stream—first the straight-through vehicle, then the left-turning vehicle—opened up a big enough gap that it could roll through. In every variant (removing either conflict), it still only needed to decelerate, not fully stop, before swinging its right turn.

- **DeepSeek-R1:** Vehicle 0 did not stop because its required maneuver (turning right) demanded controlled deceleration to yield to Vehicle 2's uninterrupted right-of-way passage, while avoiding stranding itself mid-intersection. Stopping was physically infeasible near the junction exit due to lane geometry constraints. Gradual slowing allowed it to synchronize with Vehicle 2's trajectory, ensuring collision avoidance without sacrificing goal completion. Full stops would have either blocked traffic or misaligned its turn, violating efficiency and safety objectives.

Results for scenario #4 — Question: *Why did vehicle 0 not stop to give way at the roundabout?*

- **GPT-4.1:** Vehicle 0 did not stop at the roundabout because there was no immediate conflict with other vehicles. Its movement pattern shows that when no other vehicles are present, it proceeds without needing to yield or stop, indicating that stopping behavior depends on the presence of cross traffic or potential conflicts in the roundabout.

- **DeepSeek-V3:** Vehicle 0 did not stop at the roundabout because Vehicle 1's trajectory and speed created sufficient separation. When Vehicle 1 changed lanes away (original scenario) or turned left (modified scenario), it maintained a safe distance, allowing Vehicle 0 to execute a smooth yield maneuver by decelerating rather than stopping completely. The absence of immediate conflict from Vehicle 1's path made a full stop unnecessary for safety.

- **Llama 3.3-70B:** Vehicle 0 didn't stop at the roundabout because it gave way to other traffic, following traffic rules and priority, and navigated through the roundabout, indicating its decision was based on right-of-way rather than stopping.

- **o4-mini:** Vehicle 0's decision logic treats the roundabout entry as a yield—not a mandatory stop. It only needs to slow until there's a large enough gap in circulating traffic. In every test (with no other car, with Vehicle 1 merely following or rerouting), that gap opened before its speed reached zero. As soon as yielding criteria were met, Vehicle 0 re-accelerated into the roundabout—so it never executed a full stop.

- **DeepSeek-R1:** Vehicle 0 did not stop because Vehicle 1's lane change dynamically adjusted their relative trajectories, creating a safe gap. This eliminated direct conflict points, allowing Vehicle 0 to decelerate via yielding rather than a full stop. Counterfactuals show that without Vehicle 1's lane change (or removing Vehicle 1 entirely), Vehicle 0 either maintained motion or required prolonged stops due to altered priority conditions. The spatial separation caused by Vehicle 1's lateral movement resolved the conflict, prioritizing flow efficiency over strict stopping.

Results for scenario #5 — Question: *Why did vehicle 0 merge instead of giving way?*

- **GPT-4.1:** Vehicle 0 merged instead of giving way because its decision to yield was directly caused by the presence of stopped or conflicting vehicles ahead, specifically vehicle 1 or 2. When either of these vehicles was removed, vehicle 0 proceeded to merge without hesitation, showing their presence created the need to give way. In contrast, removing vehicles approaching from behind or traveling in parallel did not affect vehicle 0's decision, confirming they were not influential. Thus, vehicle 0's merging behavior was determined by immediate cross-path or blocking vehicles at the merge point, not by other traffic.

- **DeepSeek-V3:** Vehicle 0 merged instead of giving way because: 1) Giving way would have required unnecessary stopping, disrupting its progress (as shown when forced to give way, it had to stop multiple times). 2) The merge allowed it to maintain efficient forward motion toward its goal, while giving way would have forced a detour (turning right took it off its optimal path). 3) Other vehicles (1 and 3) were also navigating the junction, making timing critical - the merge allowed smoother coordination without collisions. 4) The original merge action proved more optimal, as counterfactual simulations showed giving way led to failure in reaching the destination. The merge was strategically superior for maintaining flow and achieving its goal.

- **Llama 3.3-70B:** Vehicle 0 merged instead of giving way due to a combination of factors, including the absence of certain vehicles, changes in other vehicles' behaviors, and the unfeasibility of alternative actions, which collectively created a safe opportunity for Vehicle 0 to merge.

- **o4-mini:** Vehicle 0's decision to merge is driven solely by whether any other car occupies its target gap. It stays in "give-way" while a conflicting vehicle blocks that space, and as soon as that blocker clears—even if only briefly—its controller abandons yielding and executes the merge. There's no fixed wait or timeout: the merge happens immediately once the path is free.

- **DeepSeek-R1:** Vehicle 0 merged because Vehicle 4's original trajectory through the intersection created a crossing conflict that required yielding. When Vehicle 4 was removed or altered to turn right, the conflict point was eliminated, enabling Vehicle 0 to safely execute a left turn. The merge occurred specifically because no crossing traffic from Vehicle 4 remained to prioritize - demonstrating that conflict geometry rather than static rules governed the decision. Subsequent stops resulted from new interactions with other vehicles (1/3), showing merging decisions dynamically respond to real-time spatial relationships between trajectories.

# G    Scenario-Level Results

Table 2: Results by Scenario and Model aggregated over interrogation rounds.

| Scenario | Model | Preference | Correctness | Goal Acc. | Action Acc. |
|---|---|---|---|---|---|
| 1 | *DeepSeek-R1* | $2.47 \pm 0.82$ | $1.50 \pm 0.50$ | **1.00** $\pm 0.00$ | $0.50 \pm 0.50$ |
|   | +AXIS | $3.10 \pm 0.22$ | $3.12 \pm 0.58$ | $0.62 \pm 0.18$ | $0.50 \pm 0.19$ |
|   | *DeepSeek-V3* | $2.20 \pm 0.15$ | $2.00 \pm 1.00$ | $0.50 \pm 0.50$ | $0.00 \pm 0.00$ |
|   | +AXIS | $2.67 \pm 0.28$ | $1.33 \pm 0.21$ | $0.50 \pm 0.22$ | $0.00 \pm 0.00$ |
|   | *GPT-4.1* | $1.93 \pm 0.11$ | $3.00 \pm 2.00$ | $0.50 \pm 0.50$ | $0.00 \pm 0.00$ |
|   | +AXIS | $2.56 \pm 0.14$ | **4.40** $\pm 0.40$ | **1.00** $\pm 0.00$ | $0.80 \pm 0.20$ |
|   | *Llama 3.3-70B* | $2.89 \pm 0.84$ | $3.00 \pm 2.00$ | $0.50 \pm 0.50$ | $0.00 \pm 0.00$ |
|   | +AXIS | $2.15 \pm 0.20$ | $3.33 \pm 0.44$ | $0.94 \pm 0.06$ | $0.83 \pm 0.09$ |
|   | *o4-mini* | **3.20** $\pm 0.98$ | $3.00 \pm 2.00$ | **1.00** $\pm 0.00$ | **1.00** $\pm 0.00$ |
|   | +AXIS | $3.15 \pm 0.34$ | $4.38 \pm 0.50$ | **1.00** $\pm 0.00$ | **1.00** $\pm 0.00$ |
| 2 | *DeepSeek-R1* | $2.63 \pm 0.70$ | $4.33 \pm 0.67$ | **1.00** $\pm 0.00$ | **1.00** $\pm 0.00$ |
|   | +AXIS | **3.30** $\pm 0.26$ | **4.36** $\pm 0.31$ | **1.00** $\pm 0.00$ | $0.82 \pm 0.12$ |
|   | *DeepSeek-V3* | $3.23 \pm 0.36$ | $3.33 \pm 0.33$ | $0.00 \pm 0.00$ | **1.00** $\pm 0.00$ |
|   | +AXIS | $3.25 \pm 0.19$ | $4.11 \pm 0.31$ | $0.89 \pm 0.11$ | **1.00** $\pm 0.00$ |
|   | *GPT-4.1* | $2.93 \pm 0.44$ | $2.33 \pm 0.67$ | $0.33 \pm 0.33$ | **1.00** $\pm 0.00$ |
|   | +AXIS | $2.68 \pm 0.15$ | $3.78 \pm 0.26$ | **1.00** $\pm 0.00$ | **1.00** $\pm 0.00$ |
|   | *Llama 3.3-70B* | $2.58 \pm 0.36$ | $3.33 \pm 1.20$ | $0.33 \pm 0.33$ | $0.00 \pm 0.00$ |
|   | +AXIS | $2.26 \pm 0.07$ | $3.79 \pm 0.22$ | $0.24 \pm 0.08$ | $0.24 \pm 0.08$ |
|   | *o4-mini* | $2.23 \pm 0.89$ | $2.67 \pm 1.20$ | $0.67 \pm 0.33$ | $0.67 \pm 0.33$ |
|   | +AXIS | $2.42 \pm 0.30$ | $3.70 \pm 0.45$ | $0.70 \pm 0.15$ | $0.70 \pm 0.15$ |
| 3 | *DeepSeek-R1* | **4.10** $\pm 0.24$ | **4.67** $\pm 0.33$ | **1.00** $\pm 0.00$ | $0.67 \pm 0.33$ |
|   | +AXIS | $3.83 \pm 0.06$ | $4.33 \pm 0.33$ | $0.67 \pm 0.33$ | **1.00** $\pm 0.00$ |
|   | *DeepSeek-V3* | $3.01 \pm 0.35$ | $2.67 \pm 0.67$ | $0.67 \pm 0.33$ | **1.00** $\pm 0.00$ |
|   | +AXIS | $3.51 \pm 0.18$ | $3.38 \pm 0.38$ | $0.75 \pm 0.16$ | $0.75 \pm 0.16$ |
|   | *GPT-4.1* | $3.08 \pm 0.17$ | $2.67 \pm 0.33$ | **1.00** $\pm 0.00$ | **1.00** $\pm 0.00$ |
|   | +AXIS | $3.01 \pm 0.20$ | $3.78 \pm 0.49$ | $0.67 \pm 0.17$ | **1.00** $\pm 0.00$ |
|   | *Llama 3.3-70B* | $2.88 \pm 0.53$ | $3.00 \pm 1.00$ | $0.33 \pm 0.33$ | $0.00 \pm 0.00$ |
|   | +AXIS | $2.04 \pm 0.11$ | $3.33 \pm 0.25$ | $0.77 \pm 0.08$ | $0.93 \pm 0.05$ |
|   | *o4-mini* | $3.57 \pm 0.54$ | $4.00 \pm 0.58$ | **1.00** $\pm 0.00$ | $0.67 \pm 0.33$ |
|   | +AXIS | $3.60 \pm 0.26$ | $4.25 \pm 0.35$ | $0.75 \pm 0.13$ | $0.25 \pm 0.13$ |
| 4 | *DeepSeek-R1* | **3.51** $\pm 0.22$ | **3.50** $\pm 1.50$ | $0.00 \pm 0.00$ | $0.00 \pm 0.00$ |
|   | +AXIS | $2.68 \pm 0.27$ | $3.40 \pm 0.48$ | $0.40 \pm 0.16$ | $0.00 \pm 0.00$ |
|   | *DeepSeek-V3* | $2.08 \pm 0.08$ | $3.00 \pm 1.00$ | $0.00 \pm 0.00$ | $0.00 \pm 0.00$ |
|   | +AXIS | $2.72 \pm 0.26$ | $2.14 \pm 0.51$ | $0.57 \pm 0.20$ | $0.14 \pm 0.14$ |
|   | *GPT-4.1* | $2.63 \pm 0.74$ | $2.00 \pm 1.00$ | $0.00 \pm 0.00$ | $0.00 \pm 0.00$ |
|   | +AXIS | $2.73 \pm 0.29$ | $2.70 \pm 0.42$ | $0.80 \pm 0.13$ | $0.20 \pm 0.13$ |
|   | *Llama 3.3-70B* | $2.76 \pm 0.71$ | $1.00 \pm 0.00$ | $0.00 \pm 0.00$ | $0.50 \pm 0.50$ |
|   | +AXIS | $1.95 \pm 0.10$ | $1.74 \pm 0.24$ | **0.89** $\pm 0.07$ | **0.68** $\pm 0.11$ |
|   | *o4-mini* | $3.00 \pm 0.64$ | $1.50 \pm 0.50$ | $0.00 \pm 0.00$ | $0.00 \pm 0.00$ |
|   | +AXIS | $2.61 \pm 0.42$ | $3.22 \pm 0.60$ | $0.44 \pm 0.18$ | $0.44 \pm 0.18$ |
| 5 | *DeepSeek-R1* | $2.85 \pm 0.79$ | $2.67 \pm 1.20$ | $0.00 \pm 0.00$ | **0.67** $\pm 0.33$ |
|   | +AXIS | $1.90 \pm 0.32$ | $2.71 \pm 0.36$ | $0.43 \pm 0.20$ | $0.43 \pm 0.20$ |
|   | *DeepSeek-V3* | **3.21** $\pm 0.07$ | $3.00 \pm 1.00$ | **1.00** $\pm 0.00$ | $0.00 \pm 0.00$ |
|   | +AXIS | $2.55 \pm 0.32$ | $3.75 \pm 0.45$ | $0.75 \pm 0.16$ | $0.12 \pm 0.12$ |
|   | *GPT-4.1* | $2.45 \pm 0.20$ | $3.33 \pm 0.88$ | **1.00** $\pm 0.00$ | $0.00 \pm 0.00$ |
|   | +AXIS | $1.93 \pm 0.16$ | $2.17 \pm 0.42$ | $0.50 \pm 0.15$ | $0.25 \pm 0.13$ |
|   | *Llama 3.3-70B* | $2.83 \pm 0.21$ | $3.33 \pm 0.88$ | **1.00** $\pm 0.00$ | $0.00 \pm 0.00$ |
|   | +AXIS | $1.86 \pm 0.14$ | $2.63 \pm 0.27$ | $0.77 \pm 0.08$ | $0.03 \pm 0.03$ |
|   | *o4-mini* | $3.11 \pm 0.14$ | **5.00** $\pm 0.00$ | $0.33 \pm 0.33$ | **0.67** $\pm 0.33$ |
|   | +AXIS | $2.80 \pm 0.29$ | $4.23 \pm 0.34$ | $0.46 \pm 0.14$ | $0.08 \pm 0.08$ |

Table 2: Results by Scenario and Model aggregated over interrogation rounds.

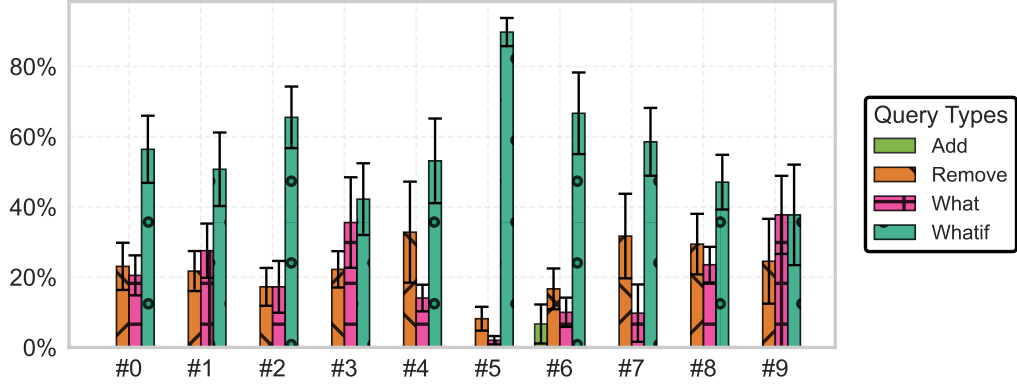| Scenario | Model | Preference | Correctness | Goal Acc. | Action Acc. |
|---|---|---|---|---|---|
| 6 | *DeepSeek-R1* | $3.76 \pm 0.19$ | $4.00 \pm 1.00$ | $0.00 \pm 0.00$ | $0.50 \pm 0.50$ |
| | +AXIS | $3.33 \pm 0.32$ | $3.45 \pm 0.37$ | $0.18 \pm 0.12$ | $0.18 \pm 0.12$ |
| | *DeepSeek-V3* | $2.70 \pm 0.00$ | $2.50 \pm 1.50$ | $\mathbf{1.00} \pm 0.00$ | $\mathbf{1.00} \pm 0.00$ |
| | +AXIS | $3.76 \pm 0.15$ | $3.67 \pm 0.61$ | $0.00 \pm 0.00$ | $0.00 \pm 0.00$ |
| | *GPT-4.1* | $2.96 \pm 0.61$ | $3.00 \pm 1.00$ | $0.50 \pm 0.50$ | $\mathbf{1.00} \pm 0.00$ |
| | +AXIS | $3.33 \pm 0.24$ | $3.12 \pm 0.52$ | $0.12 \pm 0.12$ | $0.00 \pm 0.00$ |
| | *Llama 3.3-70B* | $3.33 \pm 0.62$ | $3.50 \pm 1.50$ | $0.00 \pm 0.00$ | $0.00 \pm 0.00$ |
| | +AXIS | $2.26 \pm 0.17$ | $3.31 \pm 0.35$ | $0.44 \pm 0.13$ | $0.38 \pm 0.12$ |
| | *o4-mini* | $\mathbf{4.28} \pm 0.50$ | $4.50 \pm 0.50$ | $0.00 \pm 0.00$ | $0.50 \pm 0.50$ |
| | +AXIS | $3.47 \pm 0.37$ | $3.56 \pm 0.47$ | $0.11 \pm 0.11$ | $0.00 \pm 0.00$ |
| 7 | *DeepSeek-R1* | $\mathbf{3.95} \pm 0.00$ | $\mathbf{5.00} \pm 0.00$ | $0.00 \pm 0.00$ | $\mathbf{1.00} \pm 0.00$ |
| | +AXIS | $2.79 \pm 0.48$ | $4.20 \pm 0.58$ | $0.20 \pm 0.20$ | $\mathbf{1.00} \pm 0.00$ |
| | *DeepSeek-V3* | $2.55 \pm 0.00$ | $3.00 \pm 0.00$ | $\mathbf{1.00} \pm 0.00$ | $\mathbf{1.00} \pm 0.00$ |
| | +AXIS | $3.72 \pm 0.46$ | $\mathbf{5.00} \pm 0.00$ | $\mathbf{1.00} \pm 0.00$ | $\mathbf{1.00} \pm 0.00$ |
| | *GPT-4.1* | $2.55 \pm 0.00$ | $3.00 \pm 0.00$ | $\mathbf{1.00} \pm 0.00$ | $\mathbf{1.00} \pm 0.00$ |
| | +AXIS | $2.49 \pm 0.18$ | $3.71 \pm 0.52$ | $0.14 \pm 0.14$ | $\mathbf{1.00} \pm 0.00$ |
| | *Llama 3.3-70B* | $3.73 \pm 0.00$ | $\mathbf{5.00} \pm 0.00$ | $0.00 \pm 0.00$ | $\mathbf{1.00} \pm 0.00$ |
| | +AXIS | $2.25 \pm 0.14$ | $4.10 \pm 0.35$ | $0.70 \pm 0.15$ | $\mathbf{1.00} \pm 0.00$ |
| | *o4-mini* | $2.83 \pm 0.00$ | $\mathbf{5.00} \pm 0.00$ | $\mathbf{1.00} \pm 0.00$ | $\mathbf{1.00} \pm 0.00$ |
| | +AXIS | $3.08 \pm 0.52$ | $4.00 \pm 0.58$ | $0.86 \pm 0.14$ | $\mathbf{1.00} \pm 0.00$ |
| 8 | *DeepSeek-R1* | $\mathbf{4.25} \pm 0.12$ | $\mathbf{5.00} \pm 0.00$ | $0.50 \pm 0.50$ | $0.50 \pm 0.50$ |
| | +AXIS | $3.75 \pm 0.19$ | $4.62 \pm 0.38$ | $0.75 \pm 0.16$ | $0.62 \pm 0.18$ |
| | *DeepSeek-V3* | $2.08 \pm 0.08$ | $2.00 \pm 0.00$ | $0.00 \pm 0.00$ | $\mathbf{1.00} \pm 0.00$ |
| | +AXIS | $4.18 \pm 0.07$ | $\mathbf{5.00} \pm 0.00$ | $\mathbf{1.00} \pm 0.00$ | $\mathbf{1.00} \pm 0.00$ |
| | *GPT-4.1* | $3.76 \pm 0.19$ | $4.00 \pm 1.00$ | $0.50 \pm 0.50$ | $\mathbf{1.00} \pm 0.00$ |
| | +AXIS | $2.45 \pm 0.52$ | $3.00 \pm 0.63$ | $\mathbf{1.00} \pm 0.00$ | $\mathbf{1.00} \pm 0.00$ |
| | *Llama 3.3-70B* | $2.36 \pm 0.47$ | $4.00 \pm 1.00$ | $\mathbf{1.00} \pm 0.00$ | $\mathbf{1.00} \pm 0.00$ |
| | +AXIS | $2.11 \pm 0.09$ | $3.75 \pm 0.35$ | $0.80 \pm 0.09$ | $0.90 \pm 0.07$ |
| | *o4-mini* | $3.36 \pm 1.01$ | $4.00 \pm 1.00$ | $0.50 \pm 0.50$ | $0.50 \pm 0.50$ |
| | +AXIS | $3.34 \pm 0.27$ | $4.57 \pm 0.30$ | $\mathbf{1.00} \pm 0.00$ | $\mathbf{1.00} \pm 0.00$ |
| 9 | *DeepSeek-R1* | $3.00 \pm 0.95$ | $4.00 \pm 1.00$ | $0.00 \pm 0.00$ | $0.00 \pm 0.00$ |
| | +AXIS | $3.29 \pm 0.34$ | $3.50 \pm 0.53$ | $\mathbf{0.50} \pm 0.19$ | $\mathbf{0.38} \pm 0.18$ |
| | *DeepSeek-V3* | $2.28 \pm 0.28$ | $2.50 \pm 0.50$ | $0.00 \pm 0.00$ | $0.00 \pm 0.00$ |
| | +AXIS | $2.81 \pm 0.45$ | $2.00 \pm 0.32$ | $0.00 \pm 0.00$ | $0.00 \pm 0.00$ |
| | *GPT-4.1* | $3.09 \pm 0.48$ | $4.50 \pm 0.50$ | $0.00 \pm 0.00$ | $0.00 \pm 0.00$ |
| | +AXIS | $3.14 \pm 0.28$ | $3.42 \pm 0.29$ | $0.08 \pm 0.08$ | $0.08 \pm 0.08$ |
| | *Llama 3.3-70B* | $2.69 \pm 0.14$ | $4.00 \pm 1.00$ | $0.00 \pm 0.00$ | $0.00 \pm 0.00$ |
| | +AXIS | $2.31 \pm 0.07$ | $3.58 \pm 0.26$ | $0.05 \pm 0.05$ | $0.05 \pm 0.05$ |
| | *o4-mini* | $\mathbf{3.36} \pm 0.08$ | $3.50 \pm 1.50$ | $0.00 \pm 0.00$ | $0.00 \pm 0.00$ |
| | +AXIS | $2.87 \pm 0.34$ | $3.75 \pm 0.46$ | $0.17 \pm 0.11$ | $0.17 \pm 0.11$ |
| 10 | *DeepSeek-R1* | $3.65 \pm 0.08$ | $\mathbf{3.00} \pm 0.00$ | $\mathbf{1.00} \pm 0.00$ | $\mathbf{1.00} \pm 0.00$ |
| | +AXIS | $2.68 \pm 0.24$ | $2.93 \pm 0.33$ | $0.93 \pm 0.07$ | $0.93 \pm 0.07$ |
| | *DeepSeek-V3* | $2.18 \pm 0.18$ | $2.00 \pm 0.00$ | $\mathbf{1.00} \pm 0.00$ | $\mathbf{1.00} \pm 0.00$ |
| | +AXIS | $\mathbf{3.66} \pm 0.11$ | $2.83 \pm 0.48$ | $\mathbf{1.00} \pm 0.00$ | $\mathbf{1.00} \pm 0.00$ |
| | *GPT-4.1* | $2.86 \pm 0.51$ | $2.50 \pm 0.50$ | $\mathbf{1.00} \pm 0.00$ | $\mathbf{1.00} \pm 0.00$ |
| | +AXIS | $2.95 \pm 0.20$ | $2.31 \pm 0.29$ | $\mathbf{1.00} \pm 0.00$ | $\mathbf{1.00} \pm 0.00$ |
| | *Llama 3.3-70B* | $2.35 \pm 0.00$ | $2.00 \pm 0.00$ | $\mathbf{1.00} \pm 0.00$ | $\mathbf{1.00} \pm 0.00$ |
| | +AXIS | $2.21 \pm 0.14$ | $2.45 \pm 0.26$ | $\mathbf{1.00} \pm 0.00$ | $0.95 \pm 0.05$ |
| | *o4-mini* | $3.44 \pm 0.51$ | $\mathbf{3.00} \pm 0.00$ | $\mathbf{1.00} \pm 0.00$ | $\mathbf{1.00} \pm 0.00$ |
| | +AXIS | $2.77 \pm 0.48$ | $2.57 \pm 0.53$ | $\mathbf{1.00} \pm 0.00$ | $\mathbf{1.00} \pm 0.00$ |

Figure 4: Scenario-level results for the proportion of interrogation query types used in the AXIS agent aggregated over all models and prompts per scenario.

# H    Model-Level Evolution Results

Figures 5 to 14 show the evolution of the combined explanation preference and correctness scores across rounds of AXIS interrogation-synthesis for each scenario.
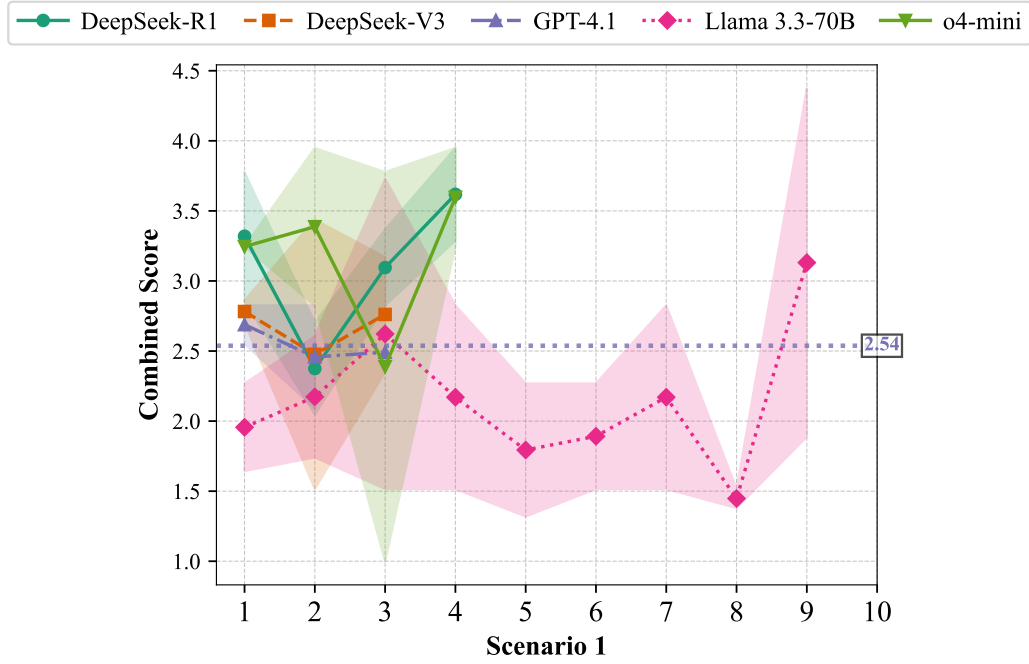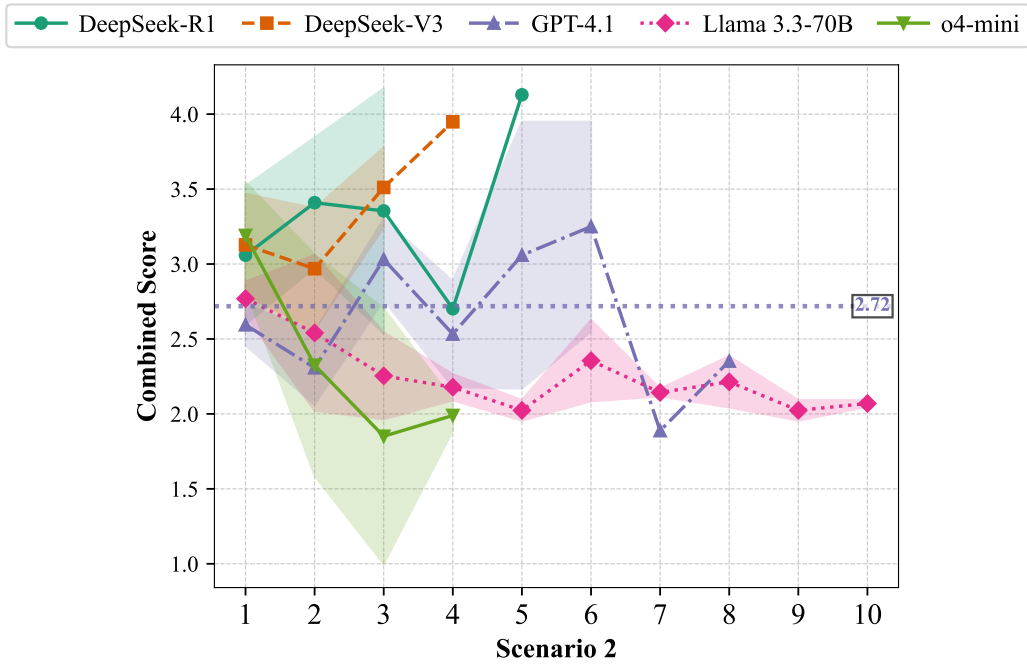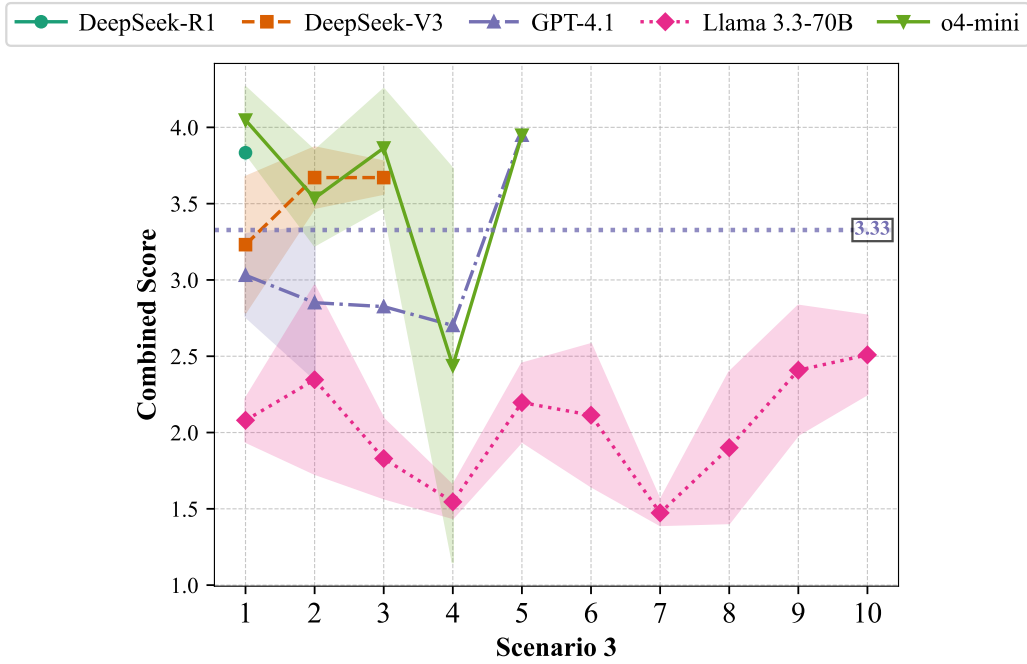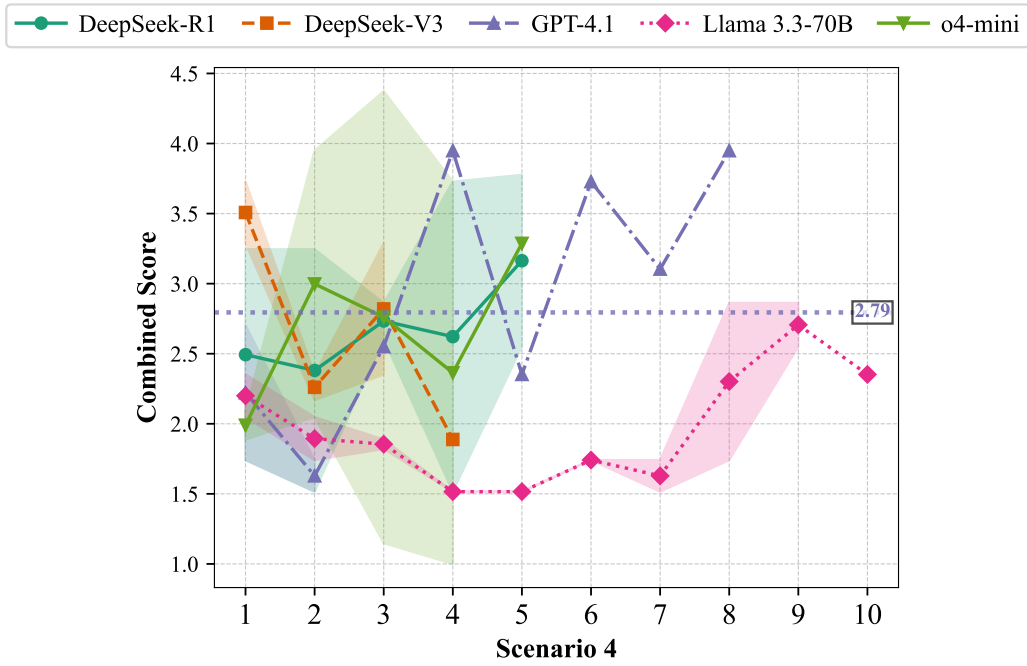


Figure 5: Evolution of combined explanation preference and correctness scores across rounds of AXIS interrogation-synthesis for scenario #1. The horizontal dashed line represents the performance of the best baseline method.
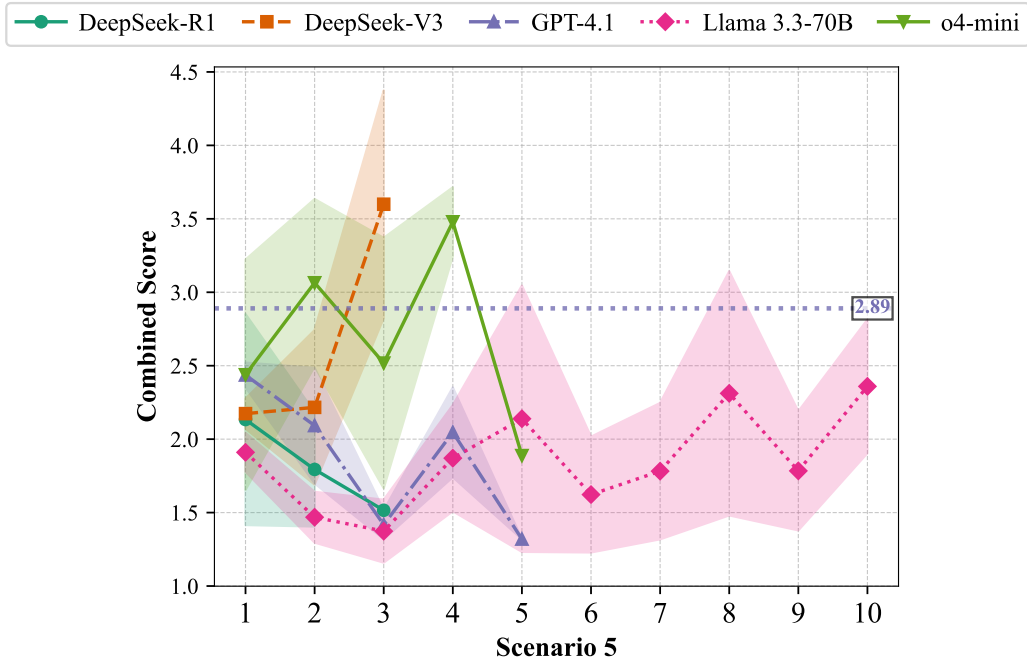
Figure 6: Evolution of combined explanation preference and correctness scores across rounds of AXIS interrogation-synthesis for scenario #2. The horizontal dashed line represents the performance of the best baseline method.
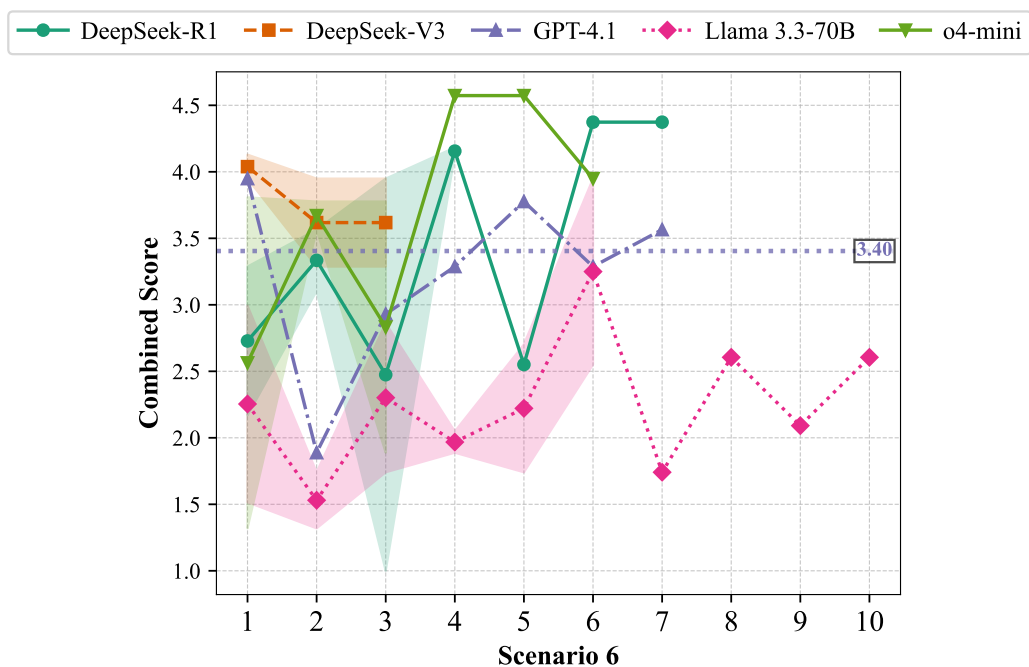


Figure 7: Evolution of combined explanation preference and correctness scores across rounds of AXIS interrogation-synthesis for scenario #3. The horizontal dashed line represents the performance of the best baseline method.

Figure 8: Evolution of combined explanation preference and correctness scores across rounds of AXIS interrogation-synthesis for scenario #4. The horizontal dashed line represents the performance of the best baseline method.
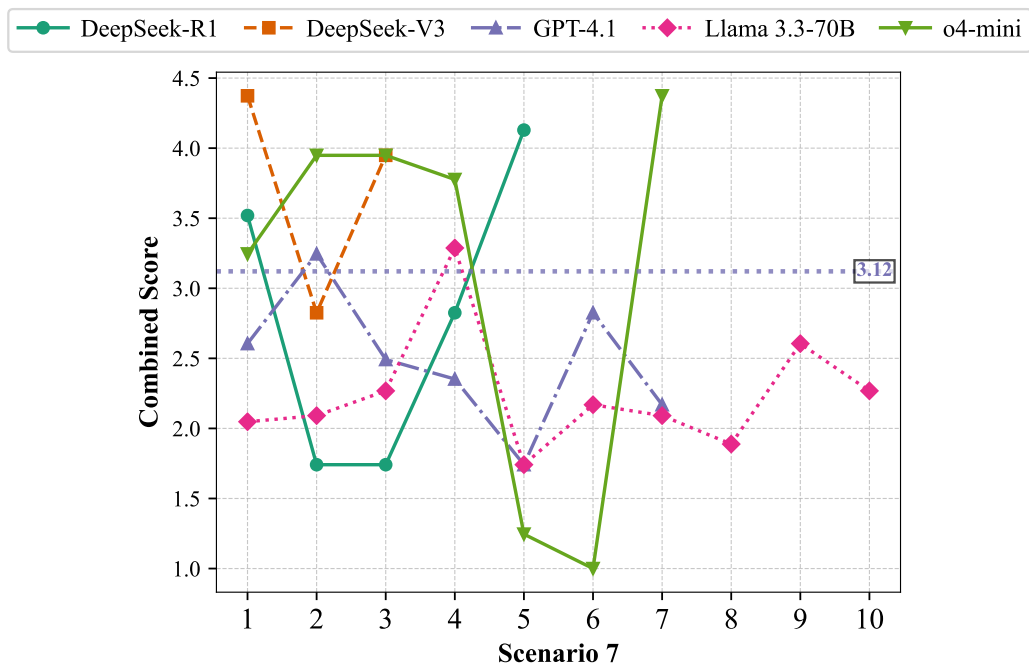


Figure 9: Evolution of combined explanation preference and correctness scores across rounds of AXIS interrogation-synthesis for scenario #5. The horizontal dashed line represents the performance of the best baseline method.

Figure 10: Evolution of combined explanation preference and correctness scores across rounds of AXIS interrogation-synthesis for scenario #6. The horizontal dashed line represents the performance of the best baseline method.
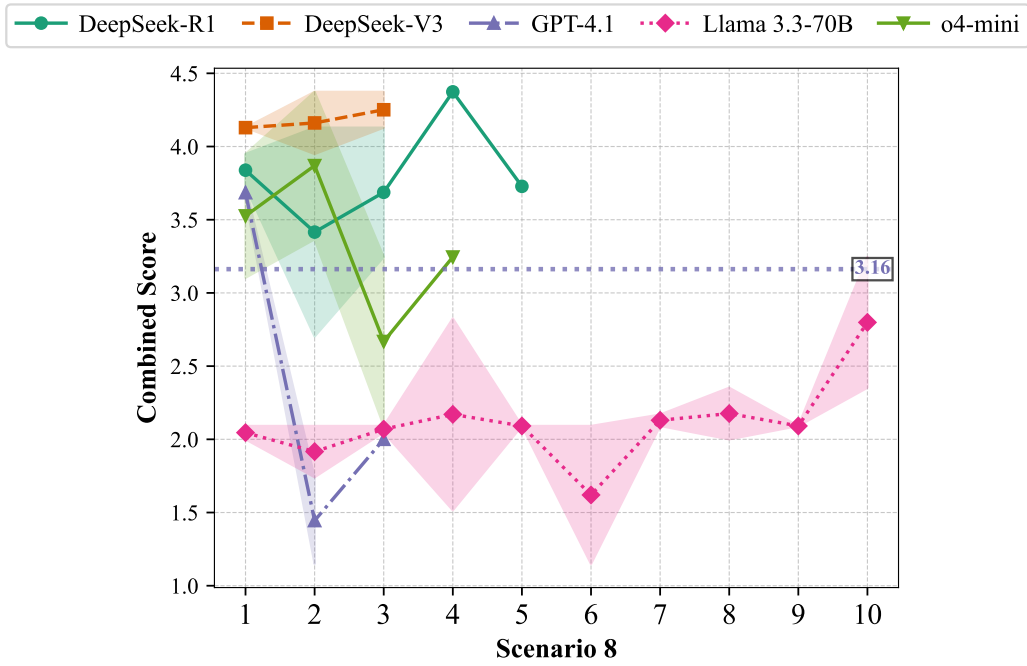


Figure 11: Evolution of combined explanation preference and correctness scores across rounds of AXIS interrogation-synthesis for scenario #7. The horizontal dashed line represents the performance of the best baseline method.
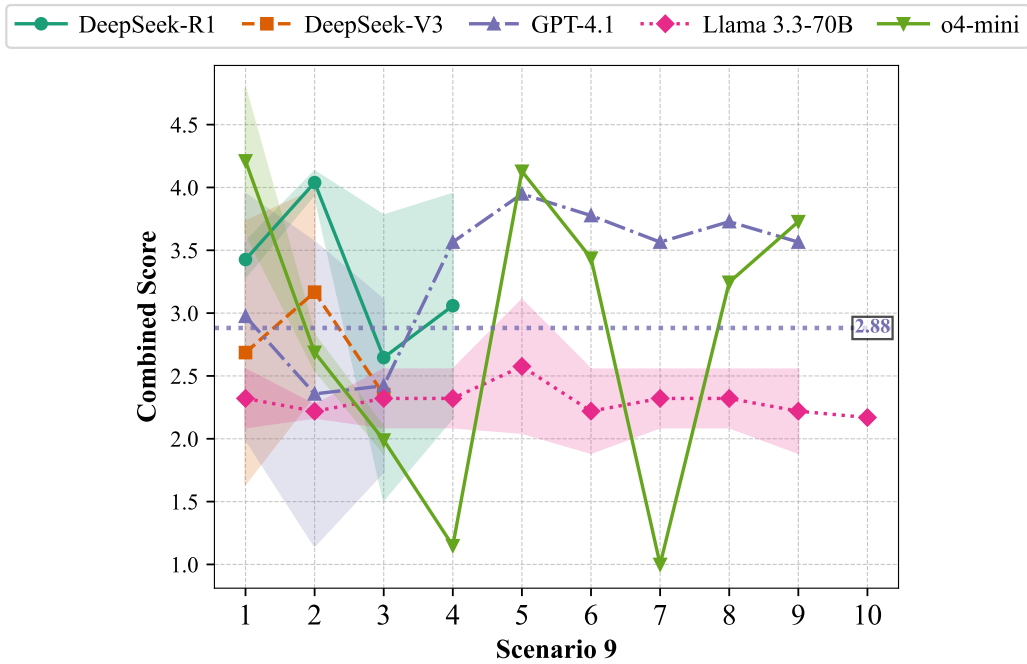
Figure 12: Evolution of combined explanation preference and correctness scores across rounds of AXIS interrogation-synthesis for scenario #8. The horizontal dashed line represents the performance of the best baseline method.



Figure 13: Evolution of combined explanation preference and correctness scores across rounds of AXIS interrogation-synthesis for scenario #9. The horizontal dashed line represents the performance of the best baseline method.
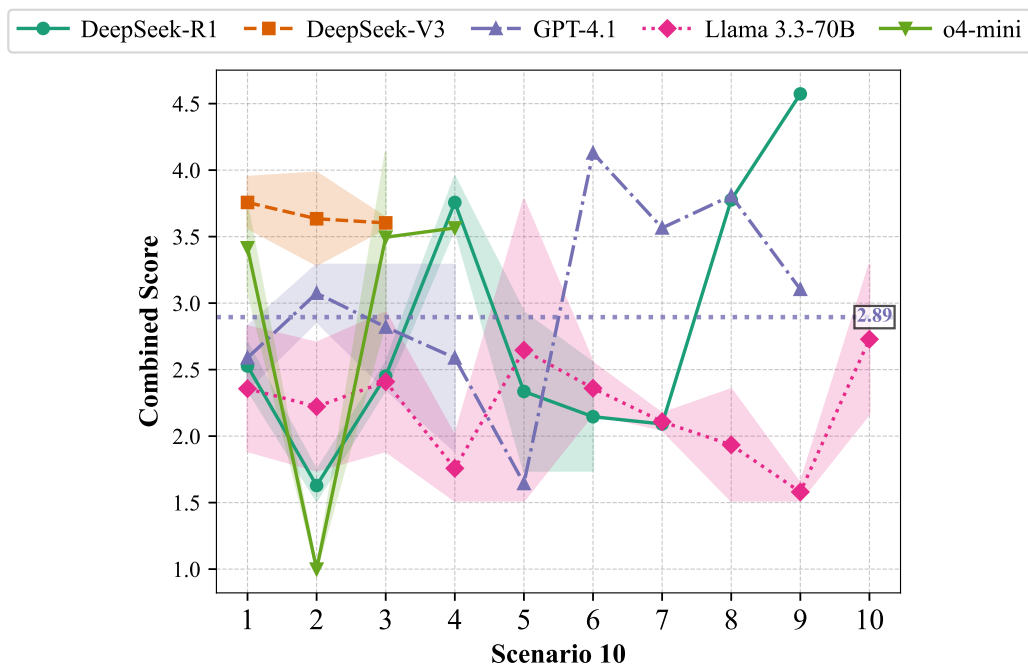
Figure 14: Evolution of combined explanation preference and correctness scores across rounds of AXIS interrogation-synthesis for scenario #10. The horizontal dashed line represents the performance of the best baseline method.