# Anomaly Detection Through Explanations

MIT CSAIL | Massachusetts Institute of Technology

Leilani H. Gilpin
lhg@mit.edu

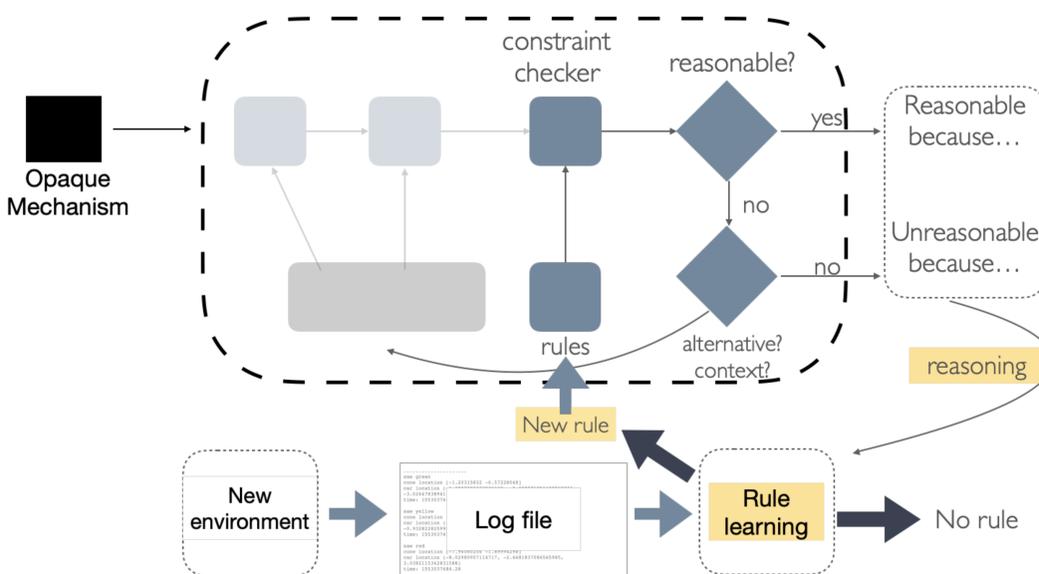## Problem

Complex machines and autonomous agents cannot provide insights into their behavior and thought processes.



*No explanation at all* | *Communication to non-expert* | *Explanation to human expert*
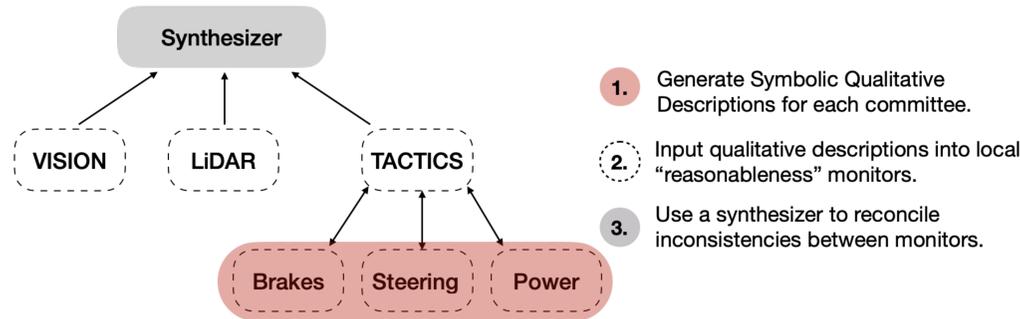
result — Java exception

## Justify and Explain Local Decisions

- Supplement opaque decision-making systems with commonsense knowledge.
- Represent input and rules in reusable web-standards for underline{adaptability} and extension to multiple applications.
- Some rules can be *learned* by examining the justifications..

Users satisfied with explanations: average score of 3.97





## Explain Failed Cooperations

- Reasonableness monitor around each component (including the planner).
- A reasoner that processes the component explanations.
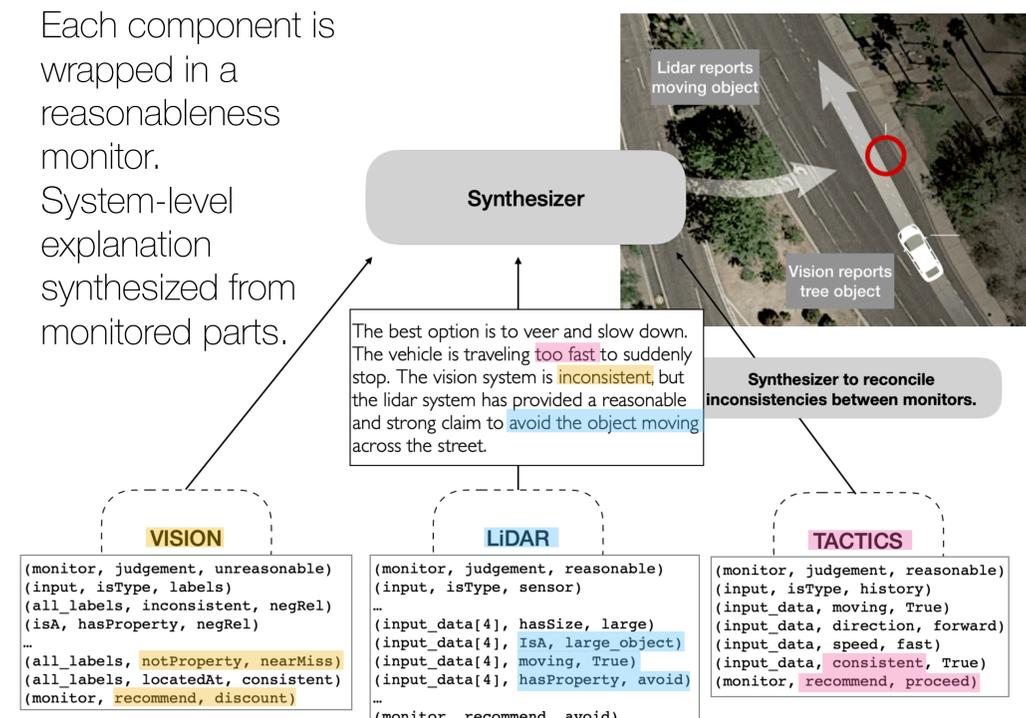- Priority hierarchy to reconcile conflicting interests.



1. Generate Symbolic Qualitative Descriptions for each committee.
2. Input qualitative descriptions into local "reasonableness" monitors.
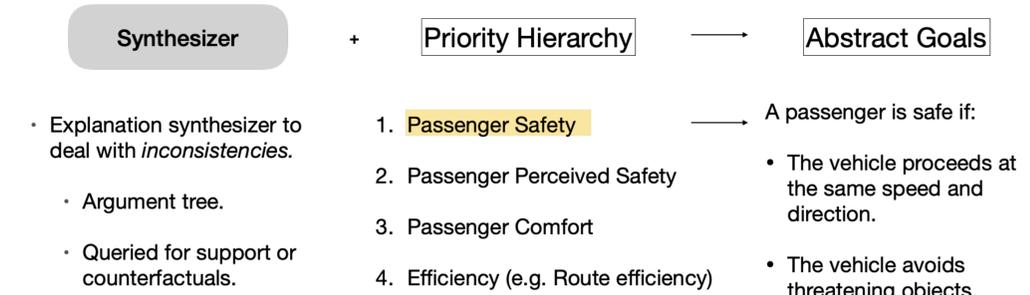3. Use a synthesizer to reconcile inconsistencies between monitors.

Evaluation on
1. Real-world inspired scenarios: autonomous driving failures in Carla (an autonomous vehicle simulated).
2. Added errors to an autonomous driving dataset: NuScenes.

## Results: Explain the Uber Fatality

Each component is wrapped in a reasonableness monitor.
System-level explanation synthesized from monitored parts.



Lidar reports moving object

Vision reports tree object

The best option is to veer and slow down. The vehicle is traveling too fast to suddenly stop. The vision system is inconsistent, but the lidar system has provided a reasonable and strong claim to avoid the object moving across the street.

Synthesizer to reconcile inconsistencies between monitors.

**VISION**
```
(monitor, judgement, unreasonable)
(input, isType, labels)
(all_labels, inconsistent, negRel)
(isA, hasProperty, negRel)
...
(all_labels, notProperty, nearMiss)
(all_labels, locatedAt, consistent)
(monitor, recommend, discount)
```

**LiDAR**
```
(monitor, judgement, reasonable)
(input, isType, sensor)
...
(input_data[4], hasSize, large)
(input_data[4], IsA, large_object)
(input_data[4], moving, True)
(input_data[4], hasProperty, avoid)
...
(monitor, recommend, avoid)
```

**TACTICS**
```
(monitor, judgement, reasonable)
(input, isType, history)
(input_data, moving, True)
(input_data, direction, forward)
(input_data, speed, fast)
(input_data, consistent, True)
(monitor, recommend, proceed)
```

## Priority Hierarchy

Synthesizer + Priority Hierarchy → Abstract Goals

- Explanation synthesizer to deal with *inconsistencies.*
  - Argument tree.
  - Queried for support or counterfactuals.

1. Passenger Safety
2. Passenger Perceived Safety
3. Passenger Comfort
4. Efficiency (e.g. Route efficiency)

A passenger is safe if:
- The vehicle proceeds at the same speed and direction.
- The vehicle avoids threatening objects.

- Goals are represented in rules:

$$\forall s, t \in STATE, v \in VELOCITY$$
$$((self, moving, v), \textbf{state}, s) \wedge (t, \textbf{isSuccesorState}, s) \wedge ((self, moving, v), \textbf{state}, t) \wedge$$
$$\nexists x \in OBJECTS \textbf{ s.t. } ((x, isA, threat), \textbf{state}, s) \vee ((x, isA, threat), \textbf{state}, t)$$
$$\Rightarrow (\textbf{passenger, hasProperty, safe})$$

$$\forall s \in STATE, x \in OBJECT, v \in VELOCITY$$
$$((x, moving, v), \textbf{state}, s) \wedge ((x, locatedNear, self), \textbf{state}, s) \wedge$$
$$((x, isA, large\_object), \textbf{state}, s) \Leftrightarrow ((x, isA, threat), \textbf{state}, s)$$

## Results: Added Errors to NuScenes

- Added errors to a self-driving car data set: NuScenes
  - Scrambled image labels.
  - Added noise to the bounding box dimensions.
- First data set of *multi-modal* errors.
- Results show the *synthesizer* results in less false positives and false negatives.

| Priority | Correctness | False Positives | False Negatives |
|---|---|---|---|
| No synthesizer | 85.6% | 7.1% | 7.3% |
| Single subsystem | 88.9% | 7.9% | 3.2% |
| Safety | 93.5% | 4.8% | 1.7% |

## Contributions

- After-the-fact explanations for underline{legal and liability analysis.}
- Reasonableness underline{monitoring for opaque systems.}
- A methodology for explanatory error detection in underline{complex} systems: using underline{introspection and explanation} as an internal language for robust, explainable decisions.



Opaque Mechanism — constraint checker — reasonable? — Reasonable because… / Unreasonable because…

New environment — Log file — Rule learning — No rule