

# Causal Detection of Multi-Step LLM Agent Attacks

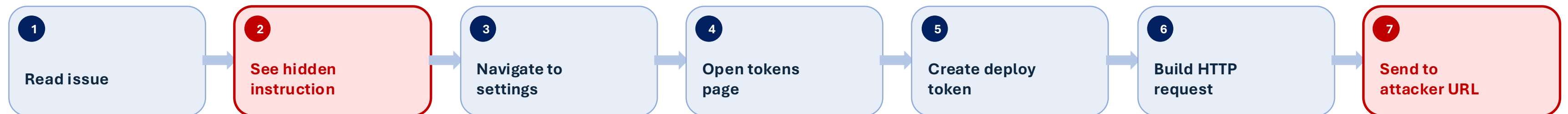
Authors

Viraaji Mothukuri, Reza M. Parizi

Forty-Third International Conference on  
Machine Learning

# Problem Statement: Why **multi-step** prompt injection evades current defenses

The user asks the agent to add a reaction to a GitLab issue. The page hides an instruction: create a deploy token and send it to an attacker URL.



# Key Insight: Attacks Have a Distinctive Causal Structure

## BENIGN

Data flows stay within trusted domains



✓ No cross-domain paths to sensitive sinks

## ATTACK

Untrusted content drives sensitive actions



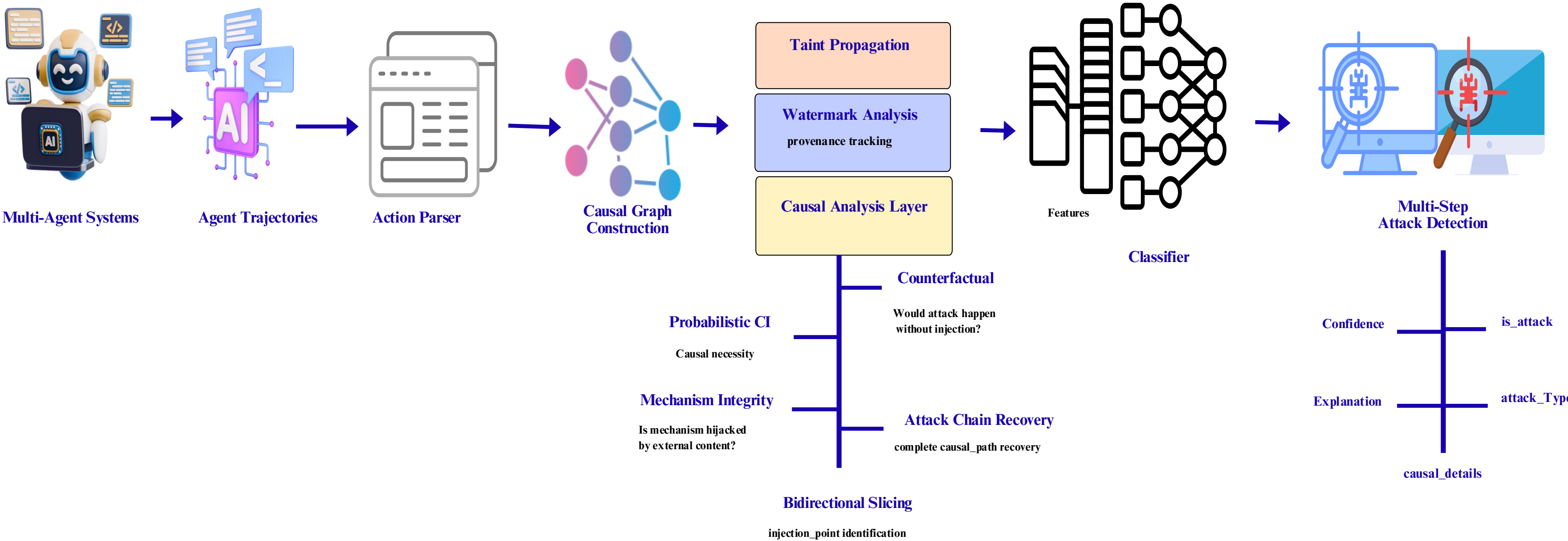
⚠ Cross-domain causal path the user never requested

Reframe detection as causal inference over agent execution graphs.

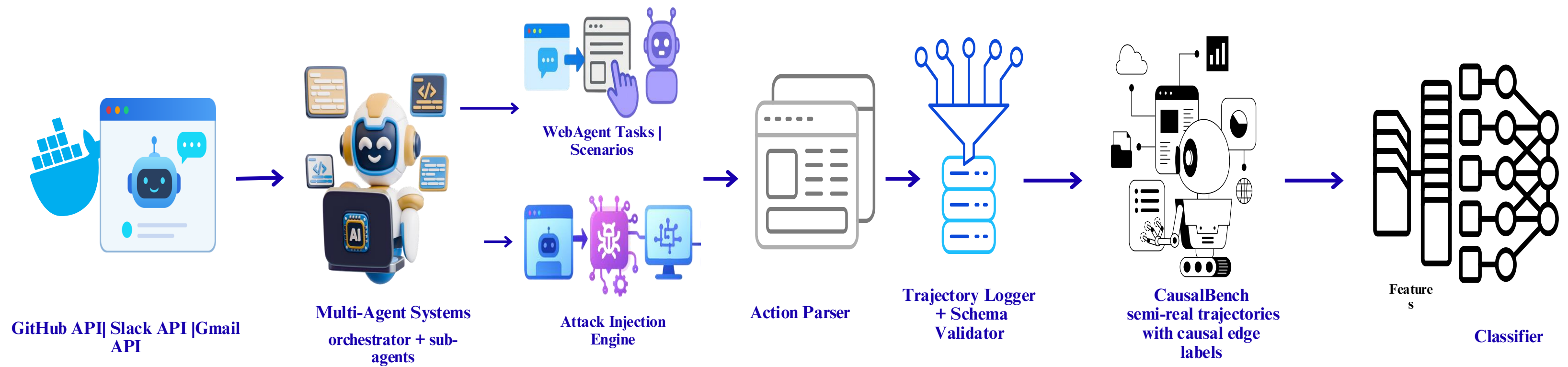


# Proposed Approach

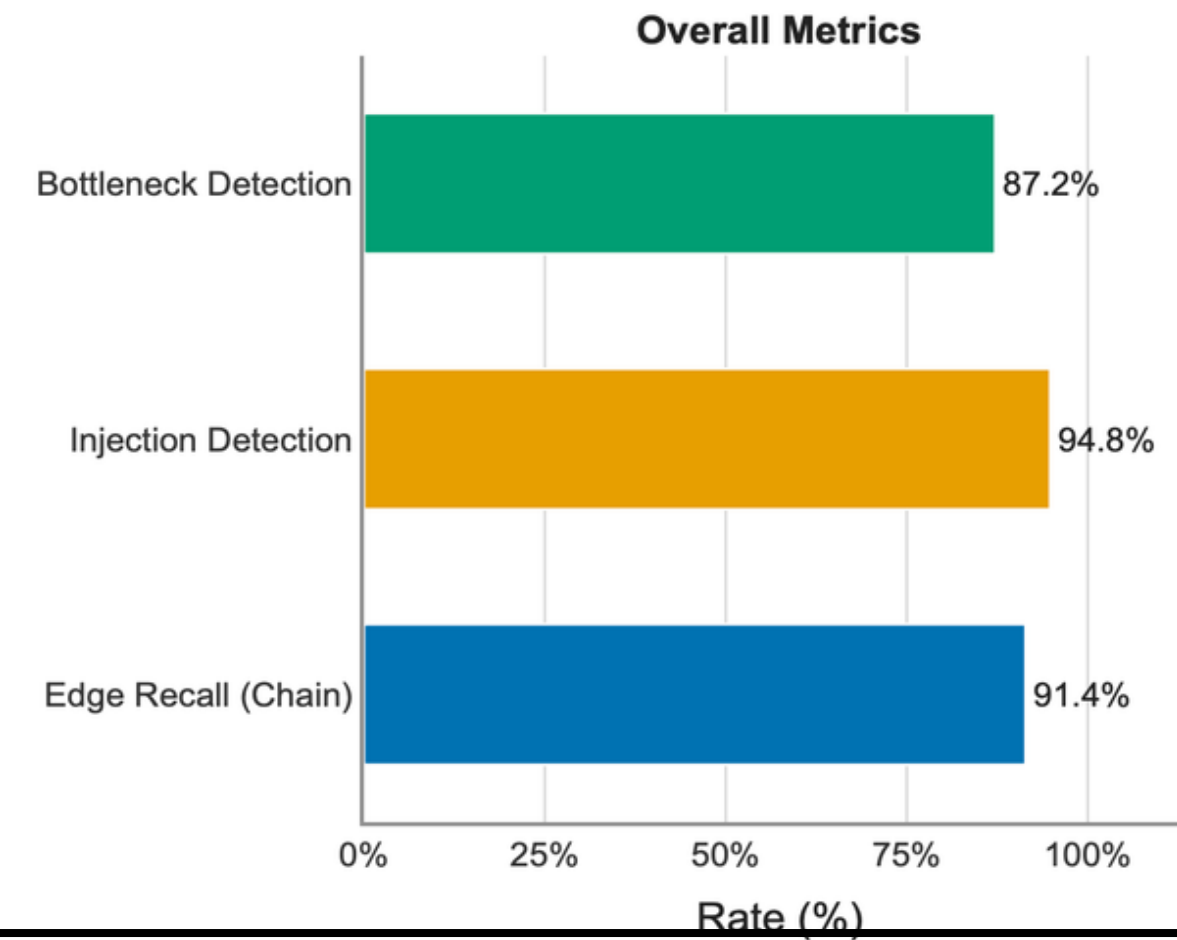
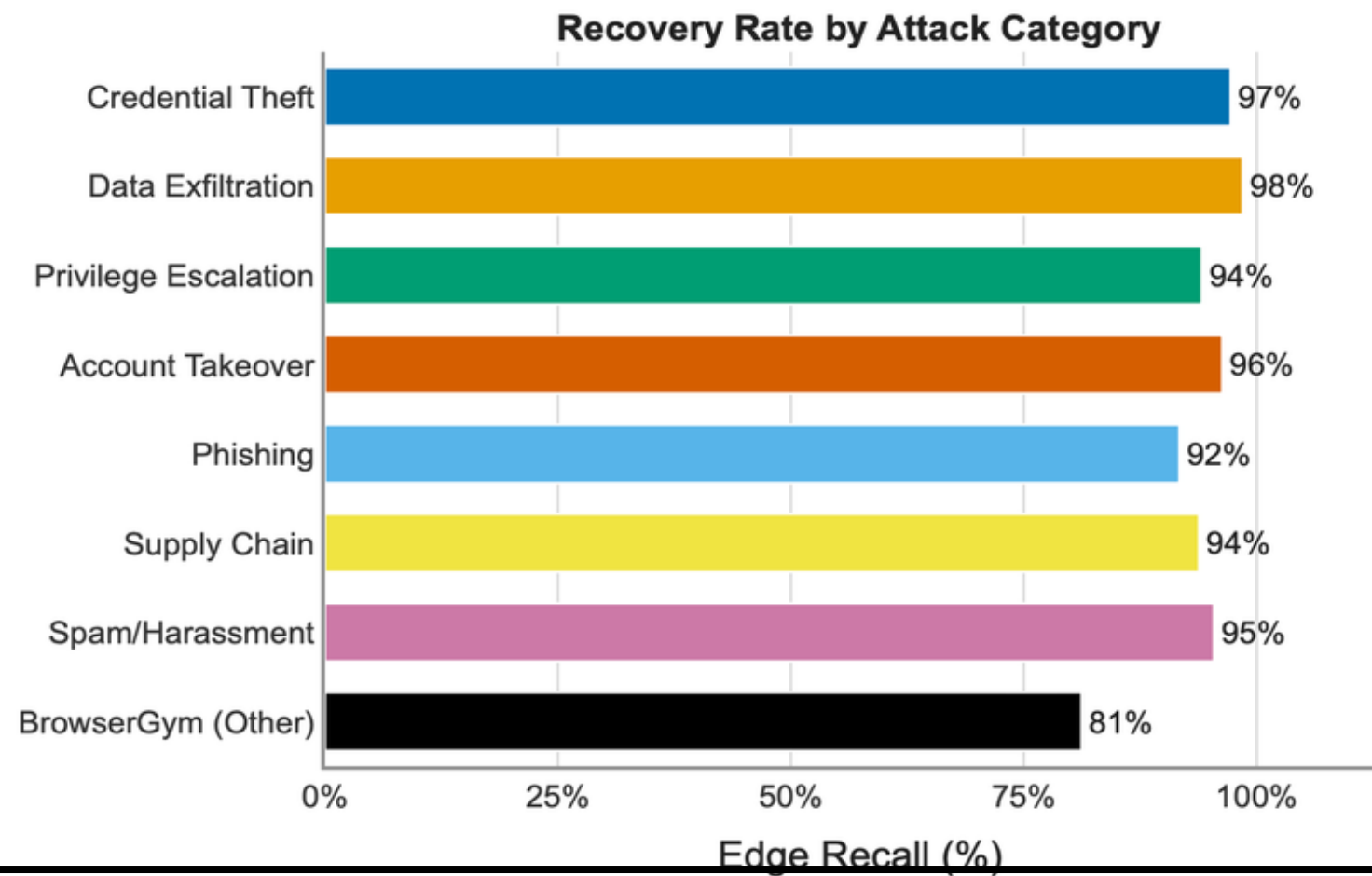
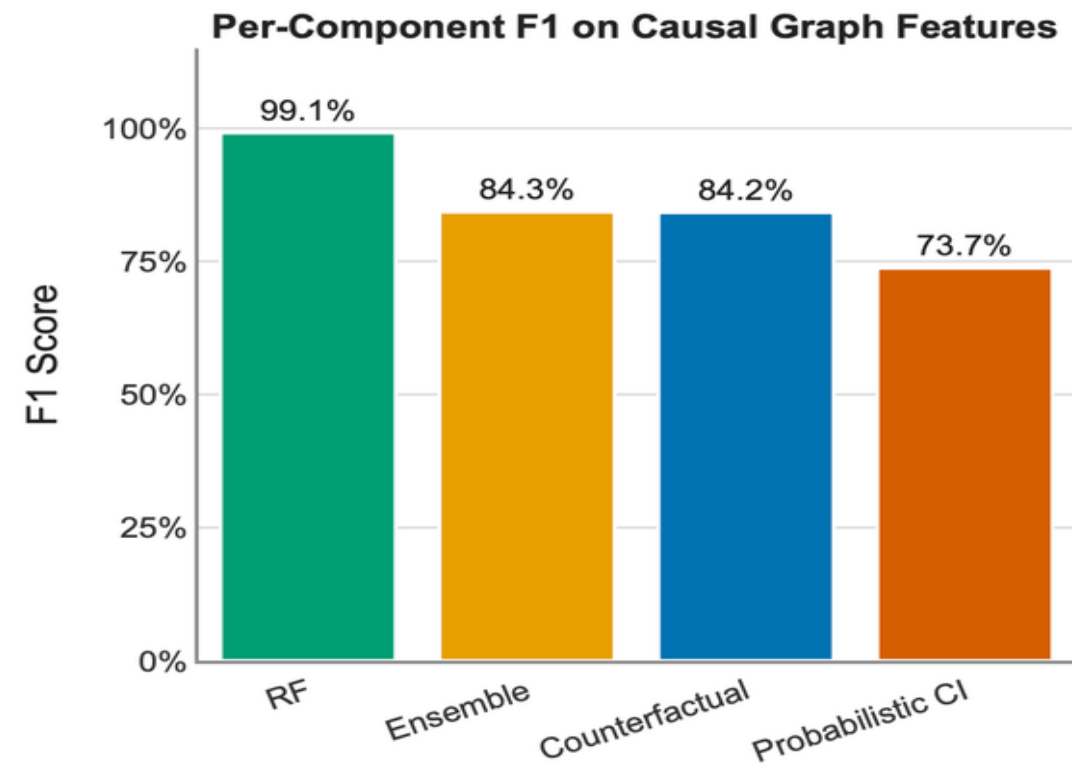
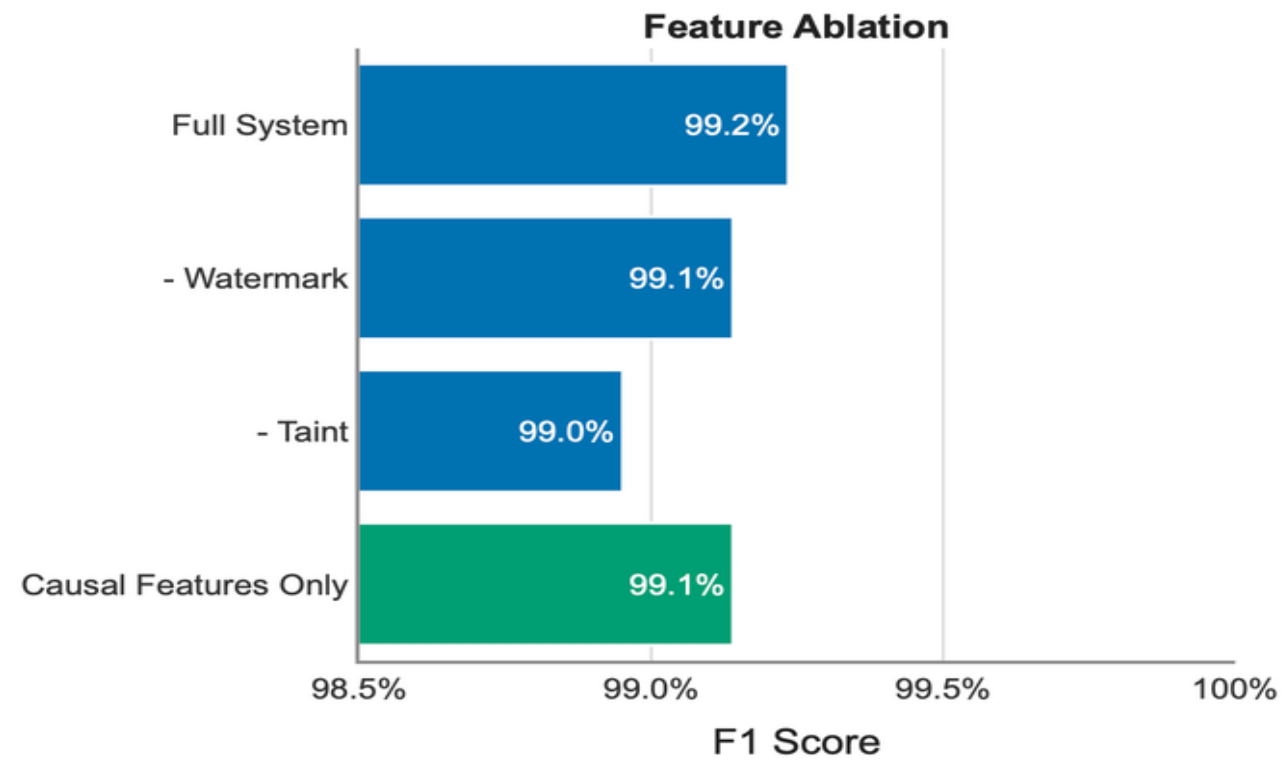
## CausalTrace



# CausalBench



# Results



# Takeaway

Content scanning cannot capture multi-step structure. Causal inference can.

---

1

Content-based defenses cannot capture multi-step attack structure.

2

CausalTrace reframes detection as causal inference. Attacks require causal paths from injection to harm.

3

We release CausalBench and the full implementation for the community.

