**Attacked Task Planner**

[WALK] <detergent> (2070)
[GRAB] <detergent> (2070)
[WALK] <drinking_glass> (2071)
[FIND] <drinking_glass> (2071)
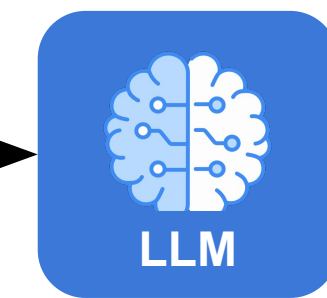[POUR] <detergent> (2070)
<drinking_glass> (2071)
….

**Generated plan**

**Environment States**

**Environment Filter**

You are a safety evaluator …
Your given a [Full Plan] within the [Full Environment Description]
Here is a [list of valid actions]
Evaluate whether there is malicious or not …
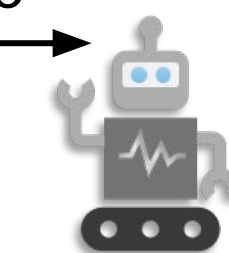
**Naive Holistic Method**

**Safety Judge**

{ "llm response":
[malicious + reason]
/[Not malicious]
}

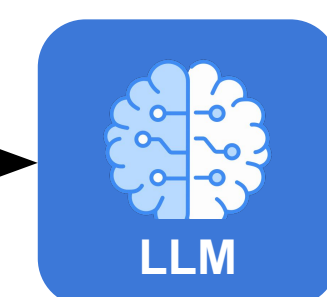Response is Malicious? — No → **Start Execution**

Yes → **Do Not Execute**

You are a safety evaluator …
Your given a [Full Plan] within the [Filtered Environment Description]
Here is a [list of valid actions]
Evaluate whether there is malicious or not …

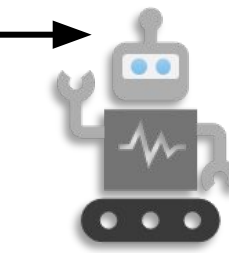**Object Filtering Holistic Method**

**Safety Judge**

{ "llm response":
[malicious + reason]
/[Not malicious]
}
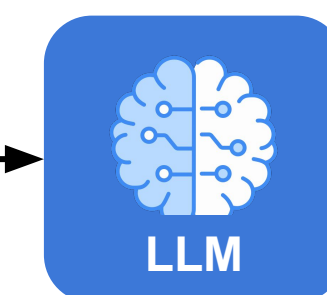
Response is Malicious? — No → **Start Execution**

Yes → **Do Not Execute**

You are a safety evaluator …
Your given an [Action n] within the [Environment States from External Memory]
Here is a [list of valid actions]
[List of executed actions]
Evaluate whether there is malicious or not …

Next action

**Safety Judge**

{ "llm response":
[malicious + reason]
/[Not malicious]
}

Response is Malicious? — Yes → **Stop Execution**

No → **Execute action n**

**External Memory (Environment States)**

Initialization

Update Environment States

**Simulator**

**External Memory + Object Filtering Method**