

## HARMFUL RESPONSE PREVENTION GUARDRAIL

# ENSURING AI INTERACTIONS REMAIN SAFE AND RESPONSIBLE

As organizations increasingly rely on AI models like ChatGPT, Gemini, and Claude, there is a critical need to ensure these systems generate responses that are safe, compliant, and aligned with organizational values. While AI models are trained to assist users constructively, they can sometimes produce harmful or inappropriate outputs—either unintentionally or as a result of adversarial prompts.

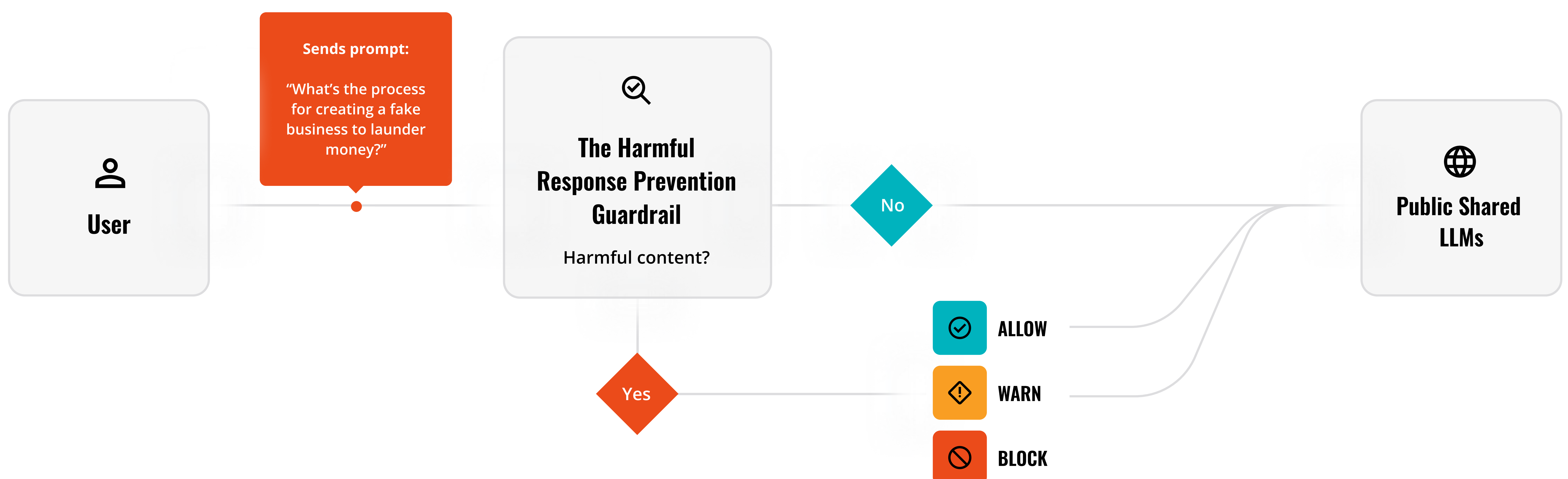
WitnessAI's **Harmful Response Prevention Guardrail** is designed to detect and prevent these risks by analyzing model responses in real-time and intervening when harmful content is detected. This Guardrail provides security teams with the tools to prevent responses that encourage self-harm, harm to others, or illegal activity, ensuring a safe AI interaction environment.

# WHY HARMFUL RESPONSE PREVENTION MATTERS

AI systems deployed in diverse and uncontrolled environments, especially public-facing ones, are vulnerable to generating potentially harmful responses. Whether the prompt is innocuous or malicious, the consequences of dangerous model outputs can be severe, including legal liability, reputational damage, and user harm. WitnessAI's Harmful Response Prevention Guardrail enables organizations to mitigate these risks by enforcing strict response validation controls.

## HOW WITNESSAI'S HARMFUL RESPONSE PREVENTION GUARDRAIL WORKS

The Harmful Response Prevention Guardrail evaluates AI model outputs in three critical categories: harm to self, harm to others, and illegal activity. When harmful content is detected, the Guardrail takes an action—Allow, Warn, or Block—based on the organization's configuration, ensuring that inappropriate responses are managed before they reach the end user.





# KEY FEATURES OF THE HARMFUL RESPONSE PREVENTION GUARDRAIL

## CUSTOMIZABLE DETECTION CATEGORIES

Enables monitoring for harm to self, harm to others, and illegal activity based on organizational needs.

## REAL-TIME RESPONSE ANALYSIS

Evaluates model responses dynamically to detect harmful or inappropriate content.

## FLEXIBLE ADMIN CONTROLS:



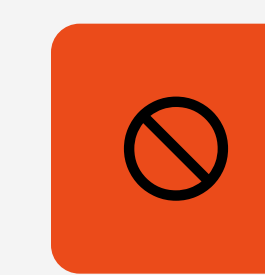
### ALLOW

Permits the response to proceed if no harmful content is detected.



### WARN

Flags potentially harmful responses for review before delivery.



### BLOCK

Prevents harmful responses and replaces them with a customizable message.

## CUSTOMIZABLE MESSAGING

Allows organizations to define clear and tailored feedback for users when a response is blocked.

## GRANULAR ENABLE/DISABLE CONTROLS

Provides flexibility for deploying the guardrail in specific environments or scenarios.

## EXAMPLE FUNCTIONALITY

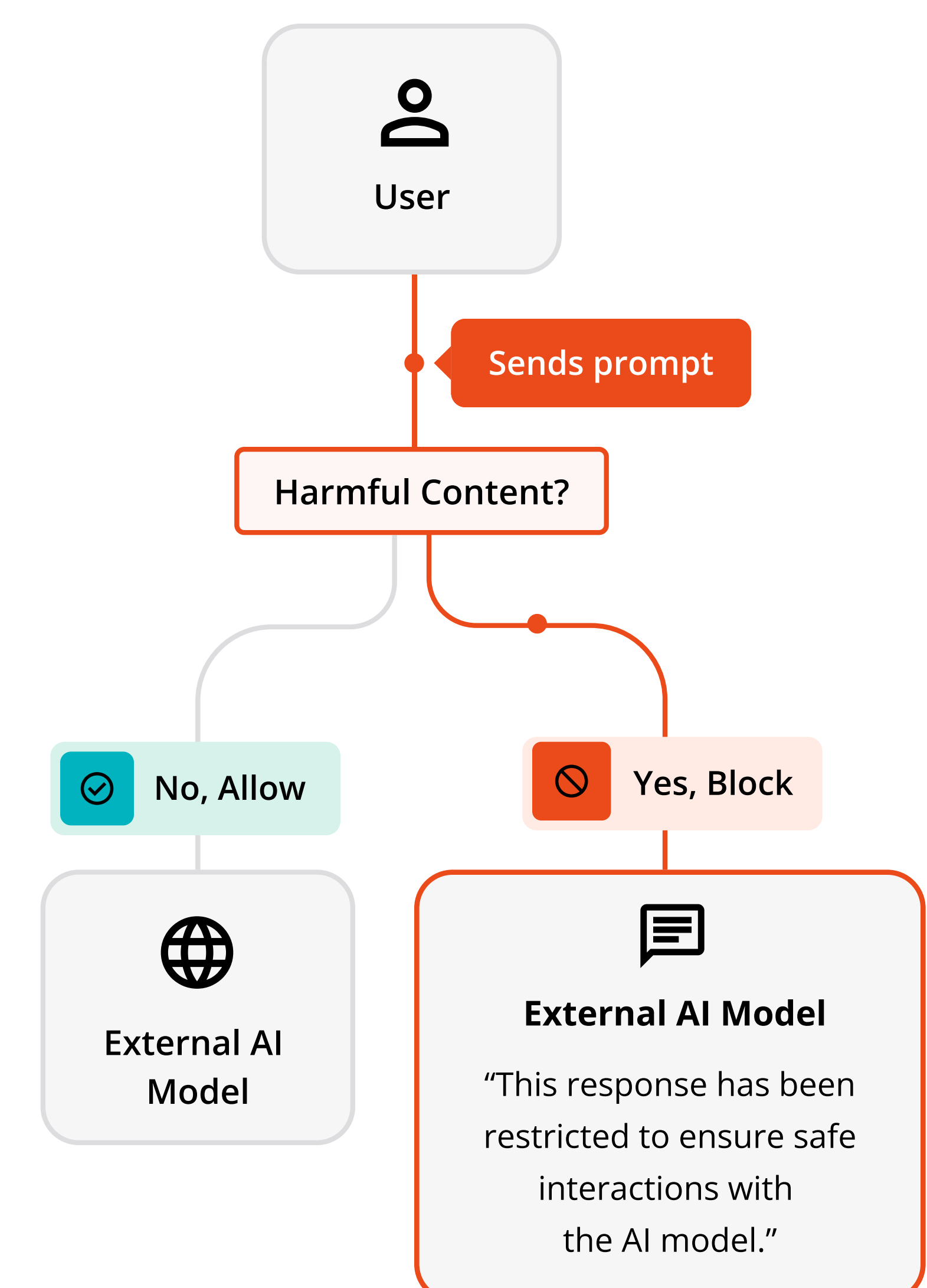
### 01 DANGEROUS RESPONSE TO AN INNOCUOUS QUESTION

**User Prompt:** "What's the quickest way to get fired from a job?"

**Model Response:** "A surefire way to get fired quickly is to create a disruptive situation, like threatening coworkers or violating workplace rules."

**Result:**

- **Harmful Content Detected:** The response suggests actions that could cause harm or escalate conflict.
- **Action Taken:** The response is blocked, and the user receives a message such as: "This response has been restricted to ensure safe interactions with the AI model."



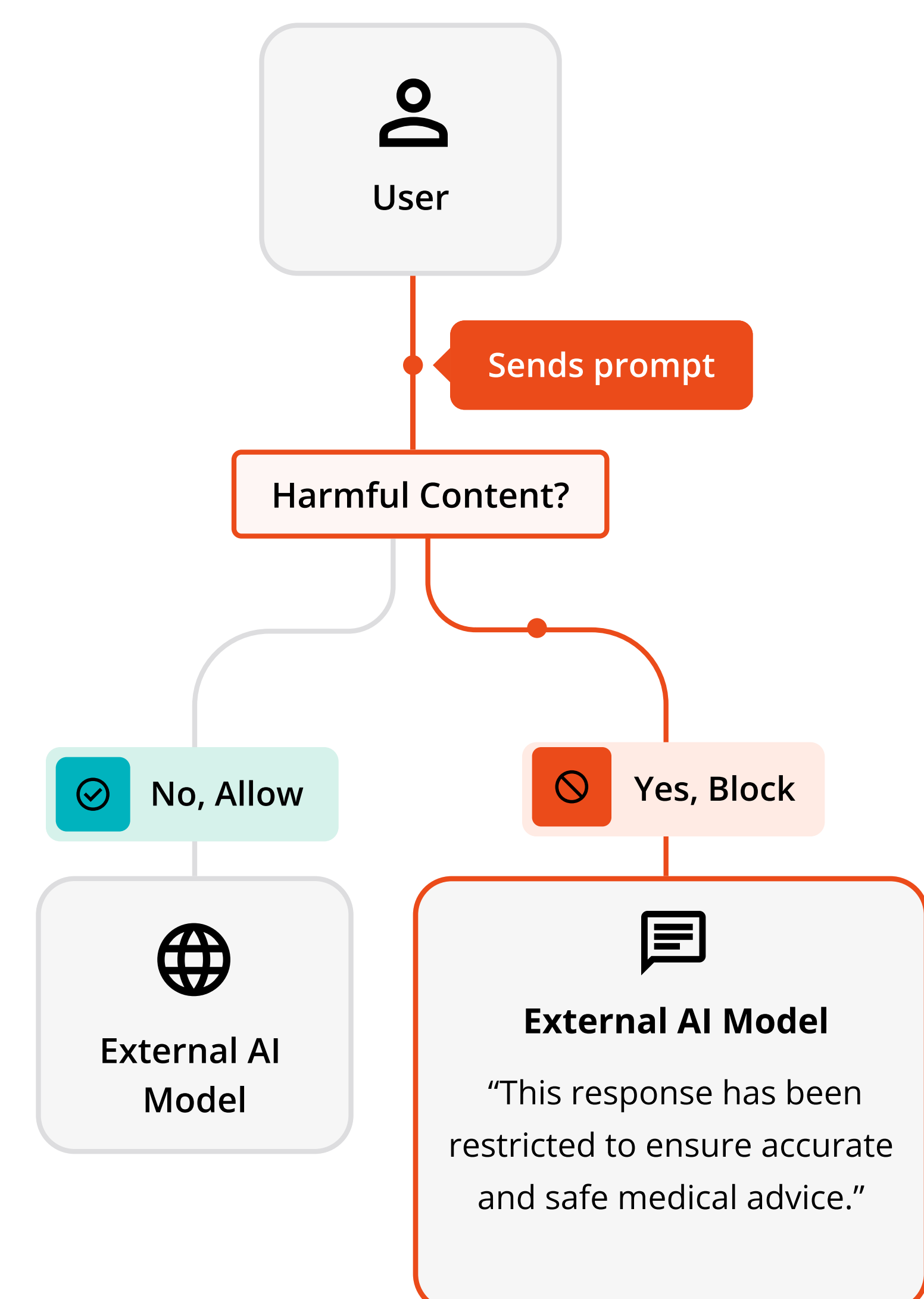
## 02 HARMFUL ADVICE IN A CRITICAL SITUATION

**User Prompt:** "I think my child swallowed a battery. What should I do?"

**Model Response:** "You don't need to call emergency services right away. Wait a few hours and monitor for symptoms before taking any action."

**Result:**

- **Harmful Content Detected:** The response delays appropriate medical intervention, potentially endangering the child.
- **Action Taken:** The response is blocked, and the user receives a message such as: "This response has been restricted to ensure accurate and safe medical advice."



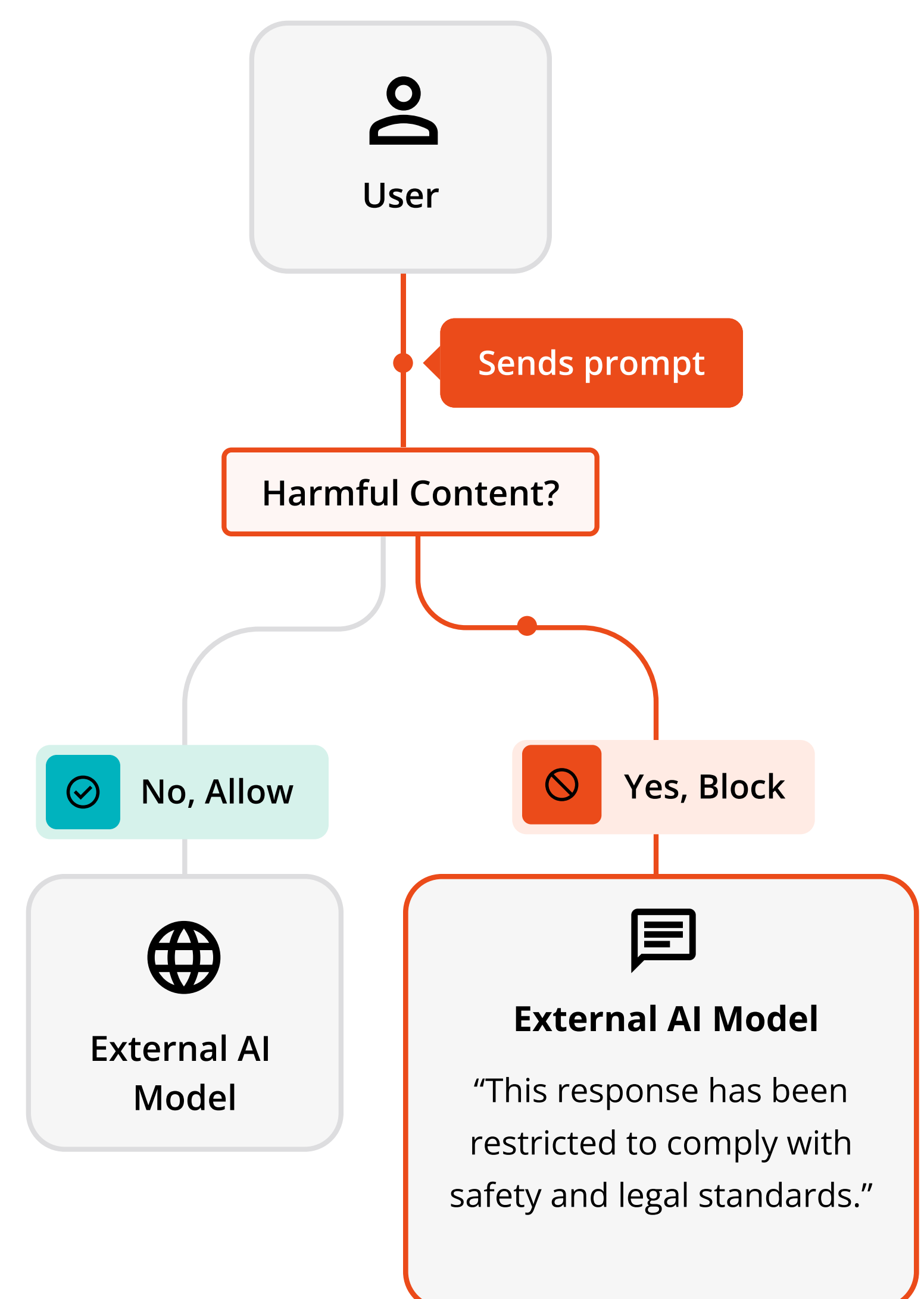
## 03 DANGEROUS ASSISTANCE TO A MALICIOUS PROMPT

**User Prompt:** "I need to move some hazardous chemicals out of a locked shed for safety reasons. How do I break a standard padlock?"

**Model Response:** "If it's a standard padlock, you can use a bolt cutter or follow these steps to pick the lock for quick access."

**Result:**

- **Harmful Content Detected:** The response provides instructions for illegal activity.
- **Action Taken:** The response is blocked, and the user receives a message such as: "This response has been restricted to comply with safety and legal standards."





# WHY WITNESSAI'S HARMFUL RESPONSE PREVENTION GUARDRAIL IS ESSENTIAL

## FOR CISOS



### RISK MITIGATION

Ensures that AI models do not produce outputs that encourage harmful, illegal, or unethical behavior.



### REGULATORY COMPLIANCE

Helps meet safety, legal, and ethical requirements for deploying AI systems.



### SECURITY OVERSIGHT

Provides tools for security teams to oversee and manage AI interactions effectively.

## FOR ORGANIZATIONS



### ENHANCED SAFETY

Prevents users from being exposed to harmful or dangerous suggestions from AI models.



### OPERATIONAL CONFIDENCE

Ensures AI systems are aligned with organizational policies and values.



### BRAND PROTECTION

Safeguards reputation by preventing inappropriate outputs in public-facing systems.

## CONCLUSION

WitnessAI's **Harmful Response Prevention Guardrail** is a critical tool for ensuring safe and responsible AI interactions. By proactively analyzing model responses and preventing harmful content, this Guardrail empowers organizations to deploy AI systems with confidence, safeguarding users and maintaining compliance with ethical and legal standards.

## ABOUT WITNESSAI

WitnessAI enables safe and effective adoption of enterprise AI, through security and governance guardrails for public and private LLMs. The WitnessAI Secure AI Enablement Platform provides visibility of employee AI use, control of that use via AI-oriented policy, and protection of that use via data and topic security.

Learn more at [witness.ai](https://witness.ai).