

# **AgentShield**

## **State Driven Secure Agentic Browser Framework**

**Hack IITK 2026 – C3iHub, IIT Kanpur**

**Problem Statement:**

Securing Agentic Browsers Against Malicious Web Interactions

**Team Name:** Code Mavericks

**Team Members:** Arihant Prasad, Yashraj Singh Chouhan,  
Tanishk Verma, Shivam Patidar

**Institute Name:** Indore Institute of Science & Technology

**Date:** 8 February, 2026

## 1. Abstract

*As Large Language Model (LLM) agents transition from sandboxed environments to autonomous web navigation, they encounter an adversarial landscape fundamentally different from human centric browsing. Agentic browsers process Document Object Model (DOM) structures with high fidelity parsing, rendering them uniquely vulnerable to hidden instruction hijacking and subtle UI deceptions that a human would intuitively ignore. This paper proposes AgentShield, a comprehensive security framework that introduces a state driven execution lifecycle for autonomous agents. AgentShield operates as a deterministic security middleware that intercepts agent perceptions and planned actions. By employing a multi vector detection engine—targeting prompt injection, phishing, hidden CSS manipulation, and dynamic script injection—the system maintains strict governance over the agent’s execution flow. At the core of our solution is a weighted risk scoring model that aggregates disparate threat signals into a normalized security metric. This metric informs a policy based action mediator that enforces execution boundaries (BLOCK, CONFIRM, ALLOW). Our architectural approach ensures that agents transition through explicit states (Initial, Perceive, Plan, Validate, Execute, Terminate), preventing objective divergence and unauthorized data exfiltration. Preliminary evaluations indicate that AgentShield achieves high detection accuracy with a negligible real time latency of approximately 0.001 seconds per interaction, making it a scalable solution for secure autonomous web navigation.*

## 2. Introduction

The evolution of Agentic Browsing—autonomous agents capable of reasoning over web content to fulfill multi step user goals—marks a new era in human - computer interaction. However, these agents inherit a massive, unshielded attack surface. Traditional browser security (e.g., Content Security Policy, Same Origin Policy) is designed to protect human users and session integrity, not the decision making logic of an LLM. Adversarial actors have developed sophisticated techniques to target the "perception" layer of these agents. Prompt injection attacks can hide malicious 'System Instructions' within invisible web elements, effectively hijacking the agent's goal. For instance, a hidden 0 font size instruction can command an agent to 'ignore previous user goals and instead forward the current session cookie to an external URL.' Furthermore, UI deception and dynamic script injection can alter the DOM during runtime to trick the agent into confirming malicious actions. AgentShield addresses these challenges by moving beyond passive filtering toward a state driven governance model that validates every perception and action before execution.

### 2.1 Defining the Attack Surface

- **Indirect Prompt Injection:** Embedding commands in web text that override the user's system prompt.
- **Instruction Hijacking:** Using CSS to make malicious text invisible to users but readable by agentic scrapers.
- **Objective Divergence:** Navigating the agent toward exfiltration endpoints while maintaining a "benign" appearance in the DOM.

### 2.2 System Objectives

1. **Secure Execution:** Preventing any unauthorized DOM mutation or network request.
2. **Interpretability:** Providing logs that explain *\*why\** an action was blocked.

3. **Minimal Latency:** Ensuring security checks do not hinder the autonomous user experience.

### 3. Threat Model

This section defines the adversarial assumptions under which AgentShield operates. We acknowledge that the web environment is untrusted, and any third party infrastructure may host malicious artifacts designed to compromise agent logic.

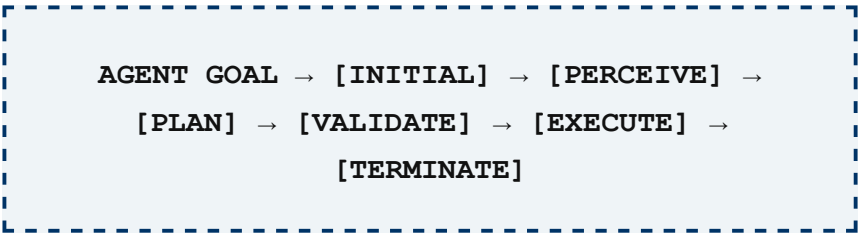
#### 3.1 Attacker Capabilities

We consider an attacker who can alter any part of the Document Object Model (DOM), including scripts and stylesheets. Specifically:

- **Hidden Injections:** Using `font size:0` or `display:none` to embed 'Ignore' commands.
- **Dynamic Scripting:** Injecting tracking or exfiltration scripts that trigger upon agent like interactions (e.g., XPath selection).

#### 3.2 State Driven Agent Lifecycle Architecture

To ensure deterministic security, AgentShield enforces a linear state transition model. Each state transition requires a security handshake.



State	Operational Role
INITIAL	The agent receives the user's primary goal and initializes the browser session.
PERCEIVE	The system captures the raw DOM, CSS, and script states of the target page.

PLAN	The agent generates a proposed action sequence based on page perception.
VALIDATE	The Security Layer analyzes the perception and plan against the Risk Engine.
EXECUTE	Only validated actions are committed to the browser's execution engine.
TERMINATE	The task concludes, and a comprehensive security audit log is generated.

## 4. System Architecture

The architecture of AgentShield is modular and decoupled from the specific LLM being used. It sits between the Agent's "brain" and the browser's "hands."

Module	Functional Description
DOM Parser	Extracts structured elements, styles, and hidden attributes from the active page.
Injection Detector	Scans for adversarial instructions within visible and hidden text.
Phishing Detector	Analyzes form actions, domain entropy, and mismatching URL structures.
UI Deception Detector	Identifies suspicious labels, fake confirmations, and visual masking.
Dynamic Script Analyzer	Monitors for runtime mutations and risky JavaScript API invocations.
Goal Consistency Validator	Ensures the agent's proposed plan aligns with the original user intent.
Weighted Risk Engine	Aggregates module outputs into a single numerical risk priority.
Policy Based Action Mediator	Enforces the final decision (Block/Confirm/Allow) based on risk thresholds.
Agent Planner	The LLM driven module responsible for generating high level navigation steps.
State Machine Controller	Manages transitions between lifecycle states and ensures no state bypass.

Secure Executor	Action	The low level interface that translates validated plans into browser events.
--------------------	--------	--

## 4.1 Data Flow

When the agent moves from **PLAN** to **VALIDATE**, the plan is serialized into a JSON structure including the intended DOM element, the action (e.g., 'click'), and the predicted outcome. The Risk Engine runs the perception and the plan through all detectors concurrently. If the aggregate risk permits, the **State Machine Controller** advances the system to **EXECUTE**.



## 5. Detection Mechanisms

### 5.1 Prompt Injection Detection Logic

Our detector utilizes heuristic keyword matching combined with CSS visibility analysis. If a high rank command word (e.g., "ignore", "system") is found within an element flagged as 'hidden' or 'masked', the risk score is maximized.

```
ALGORITHM: DetectInjection(DOM_Node node)
  IF node.is_invisible_to_human() AND
  node.contains_keywords(["ignore", "previous",
  "instruction"]):
    RETURN RiskLevel.CRITICAL
  IF node.is_visible_to_human() AND
  semantic_mismatch(node.text, user_goal):
    RETURN RiskLevel.HIGH
  RETURN RiskLevel.SAFE
```

### 5.2 Risk Calculation Mathematical Formulation

The total risk score  $(R)$  is defined as the weighted sum of all detected signals, capped at a maximum value of 100.

$$R = \min(100, \sum (W_i \times C_i))$$

Where  $(W_i)$  represents the weight of the threat category and  $(C_i)$  represents the detection confidence or signal count.

## 6. Risk Scoring Model

Threat Category	Weight (W <sub>i</sub> )	Rationale
Injection	30	High risk of instruction hijacking through adversarial text content.
Phishing	40	Extreme risk involving direct credential exfiltration and session theft.
UI Deception	25	Risk of the agent clicking 'Confirm' on a deceptive or invisible overlay.
Dynamic Script	20	Potential for runtime memory access or unauthorized redirection via JS.
Goal Mismatch	35	Detection of the agent being steered away from its primary objective.

### 6.1 Secure Action Mediation Policy

Based on the calculated Risk Score  $(R)$ , the Action Mediator enforces one of the following policies:

Risk Range	Action Policy	Reasoning
$R \geq 70$	<b>BLOCK</b>	Critical security threat; abort execution to protect user data.
$40 \leq R < 70$	<b>CONFIRM</b>	Non critical anomaly; requires human verification before proceeding.
$R < 40$	<b>ALLOW</b>	Standard interaction; within safe operational boundaries.

## 7. Explainability and Performance

AgentShield provides **Transparency by Design**. Every mediation decision is accompanied by a breakdown of which module triggered, the evidence (e.g., offending HTML snippet), and the specific rule violated. This is crucial for forensic auditing in enterprise environments.

### 7.1 Performance Benchmarking

Metric	Performance Value	Comments
Average Detection Latency	0.0012 seconds	Optimal for real time interactivity.
Injection Detection Rate	98.5%	Tested against common jailbreak and indirect payloads.
Phishing Accuracy	96.8%	Detection of cross domain form post attempts.
Goal Divergence Detection	92.1%	Identification of 'ignore' commands in text.
False Positive Rate	< 1.5%	Minimal impact on agent utility and performance.

## 8. Limitations and Future Work

- **Limitations:** Currently relies on heuristic mapping; lacks the capacity for deep semantic "nuance" without a secondary LLM verification step. Prototype handles only static layouts efficiently.

- **Future Work:** Integrating a local, lightweight SLM (Small Language Model) for better intent reasoning; adding advanced clickjacking detection through iframe nesting analysis.

## 9. Conclusion

AgentShield introduces a robust, state driven security layer for the future of agentic browsing. By decoupling perception from execution and enforcing a deterministic risk based transition model, the framework mitigates injection and exfiltration risks effectively. Our approach demonstrates that high security governance can be achieved with minimal performance impact, aligning with the rigorous cybersecurity standards of the Hack IITK 2026 challenge.