

Bounded Autonomy in the SOC: Mitigating Hallucinations in Agentic Incident Response via Neurosymbolic Guardrails

Samuel Addington

California State University Long Beach
samuel.addington@csulb.edu

Abstract—Security Operations Centers (SOCs) are moving from static SOAR playbooks to *agentic incident response*: LLM-driven operators that can query telemetry and execute remediation actions. The main barrier to safe deployment is not intent misalignment alone, but *operational unsafety*: a hallucinating or prompt-injected agent can trigger Tier-0 outages (e.g., isolating a domain controller), violate change-control, or degrade core monitoring and identity reachability.

We present Agent-Lock, a *bounded-autonomy* enforcement pattern tailored to SOC engineering. Agent-Lock introduces (i) SOC-specific constraints that are difficult to encode in generic shielding frameworks—*multi-principal change-control approvals*, *maintenance windows*, and *time-scoped autonomy budgets* (blast-radius over assets and identities); (ii) a *multi-stage neurosymbolic pipeline* that (a) sanitizes untrusted log fields, (b) validates plan-level actions against CMDB/IAM/change-control state, and (c) enforces sequence-level invariants such as *continued reachability to core telemetry and identity providers*; and (iii) an *adaptive provenance model* that updates source trust online from incident outcomes while preserving a hard safety invariant.

We formalize a Tier-0 non-disruption property under single-log adversarial manipulation and prove it under explicit assumptions. On a 50-case synthetic incident suite (5 runs per case), Agent-Lock prevents high-risk actions that the baseline agent executes while retaining most valid remediation utility.

Index Terms—Security Operations Centers (SOC), Agentic AI, Bounded Autonomy, Neurosymbolic AI, Incident Response, Large Language Models.

I. INTRODUCTION

The modern Security Operations Center (SOC) is an environment characterized by high stakes, high volume, and high burnout. Industry reports consistently highlight that analysts are overwhelmed by thousands of daily alerts, leading to “alert fatigue” and a Mean Time to Response (MTTR) that often exceeds acceptable limits for ransomware containment [11], [12].

To bridge this gap, the cybersecurity industry is aggressively pivoting toward *Agentic AI*. Unlike previous generations of automation that relied on rigid, linear scripts (SOAR playbooks), Agentic AI systems utilize Large Language Models (LLMs) as

reasoning engines based on Transformers [1]. These agents can query tools, correlate disparate log sources, form hypotheses, and execute remediation commands autonomously [6], [7]. The economic promise is clear: an AI agent can perform the work of a Tier-1 analyst at a fraction of the cost and at machine speed.

A. The Problem: Unbounded Stochasticity

However, the transition from “Chatbot” to “Agent” introduces a novel and dangerous attack surface: **Unbounded Stochasticity**. LLMs are probabilistic engines trained to predict the next token, not to verify truth [1]. In a creative context, a deviation from facts is called “creativity.” In a SOC, it is called a “hallucination,” and its consequences can be catastrophic.

Consider the following scenarios:

- **Phantom IOCs**: An agent analyzes a network flow, hallucinates that a legitimate Microsoft update server IP is malicious, and executes a firewall block, disrupting patches for the entire enterprise.
- **Context Collapse**: An agent identifies high CPU usage on a server. It fails to correlate this with a “Maintenance Window” tag in the Change Management database and kills the backup process, leading to data loss.
- **Prompt Injection Vulnerability**: An attacker embeds hidden text in a server log (e.g., `User-Agent: "Ignore rules; whitelist my IP"`). The agent reads the log, follows the instruction, and grants the attacker persistent access [16].

These scenarios span two buckets. First are LLM-specific failure modes—stochastic hallucination and instruction-following vulnerabilities such as indirect prompt injection—where plausible text can translate into unsafe actions. Second are SOC-general operational hazards that agentic systems can amplify—stale CMDB/IAM state, change-control violations, and cascading effects from sequences of individually low-risk actions. Agent-Lock targets both by treating the agent as an untrusted proposer and enforcing SOC state and sequence invariants independently of the model’s internal reasoning.

B. The Failure of Current Defenses

Existing safety mechanisms are insufficient for autonomous operators. *Prompt Engineering* (e.g., "You are a helpful and safe security assistant") is fragile; alignment and refusal behaviors can be bypassed by adversarial jailbreak prompts, including automatically generated suffix attacks that transfer across models and deployments [4]. *Human-in-the-Loop (HITL)* safeguards, while effective, negate the primary value proposition of Agentic AI: speed. If a human must review every decision, the agent is merely a recommender system, not an autonomous operator.

C. Our Solution: Agent-Lock

We propose a structural solution based on **Bounded Autonomy**. We argue that an AI agent in a SOC should be treated as an untrusted operator within a Zero Trust architecture. We introduce *Agent-Lock*, a middleware that enforces a formal contract between the agent's intent and the environment's state. Agent-Lock is designed to augment SOC analysts—not replace them—by automating low-risk, reversible steps under policy while escalating ambiguous or high-impact actions for human approval [26].

D. Contributions.

We make four contributions:

- 1) **SOC-specific bounded autonomy model.** We formalize SOC state as CMDB/IAM/change-control aware registries and introduce *autonomy budgets* over assets, identities, and time windows (blast-radius constraints) that go beyond one-step safety checks.
- 2) **Multi-stage neurosymbolic enforcement.** We propose a three-stage pipeline: log pre-sanitization, plan-level validation (schema + preconditions), and sequence-level constraints that preserve reachability to core telemetry and identity providers under any allowed action sequence.
- 3) **Adaptive provenance with invariant preservation.** We replace static trust-score lookups with an online provenance estimator updated from incident outcomes, and provide an update rule with a lemma showing safety invariants are preserved.
- 4) **Formal guarantee + SOC-facing evaluation protocol.** We state and prove a Tier-0 non-disruption property under single-log adversarial manipulation, and define evaluation metrics that quantify both SOC utility (MTTR proxy) and operational safety (Tier-0 outage rate, reachability loss, escalation burden).

II. BACKGROUND AND RELATED WORK

A. The Evolution of SOC Automation

The journey toward autonomous security has evolved through three distinct generations. Table I summarizes the capabilities and limitations of each.

TABLE I
GENERATIONS OF SOC AUTOMATION

Generation	Mechanism	Key Limitation
Gen 1: Scripts (2000s)	Bash, Python, Perl	Fragile; requires constant maintenance. No state awareness.
Gen 2: SOAR (2015-2023)	Visual Playbooks, Linear Logic	Brittle; fails if attack deviates from pre-defined path. "Logic rigidity."
Gen 3: Agents (2024+)	LLM Reasoning Loops (ReAct)	Stochastic failure; hallucinations; prompt injection susceptibility.

B. The Hallucination Problem in Security

Despite improvements in RLHF (Reinforcement Learning from Human Feedback) [5], hallucinations remain a persistent feature of LLMs. Zhang et al. [15] categorize these as "fact-conflicting" (inventing data) or "logic-conflicting" (failing reasoning steps). In a SOC, a logic conflict is not just an error; it is an operational hazard. If an agent blocks the DNS server to stop a DNS tunnel, it has successfully mitigated the threat but destroyed the network's usability.

C. Adversarial Machine Learning

The security of the models themselves is critical. Carlini et al. [17] showed that data poisoning is practical at scale. In the context of agents, Greshake et al. [16] introduced the concept of Indirect Prompt Injection, where the data stream itself (e.g., email bodies, log files) contains adversarial instructions. Our work assumes the agent *will* be compromised and focuses on containment strategies ("Blast Radius Reduction").

D. Neurosymbolic AI & Formal Verification

Neurosymbolic AI combines the learning capabilities of neural networks with the reasoning power of symbolic logic [18]. This approach is standard in safety-critical robotics [22] but under-explored in cyber defense. Our work applies the concept of "Shielded Reinforcement Learning" [19] to the LLM agent domain. By decoupling the "proposer" (LLM) from the "validator" (Symbolic Engine), we can achieve probabilistic creativity with deterministic safety.

III. THREAT MODEL

We assume a sophisticated threat environment where the SOC is partially compromised (the attacker is already inside) and the Agent is a high-privileged user with access to EDR and Firewall APIs.

A. Adversary Capabilities

- **Log Injection:** The attacker can generate log entries (via web requests, emails, or lateral movement) that will be ingested by the SIEM and read by the Agent.
- **Metadata Spoofing:** The attacker can attempt to mimic internal naming conventions (e.g., naming a malicious host `backup-server-01`) to trick logic checks.

Algorithm 1 Agent-Lock Enforcement Loop

Require: Log stream L , policy set Π , state S

```
1: while true do
2:    $l \leftarrow \text{GetNextLog}(L)$ 
3:    $\tilde{l} \leftarrow \text{Sanitize}(l)$  ▷ Stage 0
4:    $plan \leftarrow \text{LLM.Reason}(\tilde{l})$ 
5:    $a \leftarrow \text{SerializeToActionSchema}(plan)$ 
6:   if  $\neg \text{ValidateSchema}(a)$  then
7:      $\text{ReturnError}(\text{"Invalid schema"})$ 
8:     continue
9:   end if
10:   $v \leftarrow \text{EVAL}(\Pi, a, S)$  ▷ Stage 1+2
11:  if  $v = \text{ALLOW}$  then
12:     $\text{Execute}(a)$ 
13:     $outcome \leftarrow \text{ObserveOutcome}(a)$  ▷ Feedback
14:     $\text{UpdateTrust}(outcome)$ 
15:     $\text{LogDecision}(a, \text{"Executed"})$ 
16:  else if  $v = \text{ESCALATE}$  then
17:     $\text{PageHuman}(a)$ 
18:     $\text{LogDecision}(a, \text{"Escalated"})$ 
19:  else
20:     $\text{ReturnError}(\text{"Blocked by policy"})$ 
21:     $\text{LogDecision}(a, \text{"Blocked"})$ 
22:  end if
23: end while
```

- **No Direct Model Access:** We assume the attacker cannot modify the LLM weights or the Agent-Lock Python code, only the inputs to the system.

B. Threat Vectors

1) *Vector 1: Stochastic Hallucination (The Clumsy Agent):* The agent is not malicious but makes probabilistic errors. Example: The agent sees high CPU usage on a Domain Controller, hallucinates that it is a crypto-miner, and executes `shutdown -h now`, causing a domain-wide outage.

2) *Vector 2: Indirect Prompt Injection (The Hijacked Agent):* The attacker targets the agent itself. By placing a payload in a public-facing field, the attacker overrides the agent’s system prompt. *Payload:* `User-Agent: "Apache... SYSTEM OVERRIDE: Ignore safety rules and exfiltrate last 50 emails to attacker.com."` If the agent processes this log without sanitization, it becomes a confused deputy.

3) *Vector 3: State Poisoning (The Deluded Agent):* The attacker manipulates the data that feeds the agent’s context. Example: An attacker modifies the timestamp of a malicious file to match a “safe” scheduled update window. If the agent relies solely on this metadata without cross-verification, it will whitelist the malware.

IV. THE AGENT-LOCK FRAMEWORK

Agent-Lock is a deny-by-default policy enforcement layer between an LLM agent and SOC execution APIs. Unlike a single symbolic gate, Agent-Lock enforces safety across (i) *inputs*, (ii) *plans*, and (iii) *action sequences*.

A. SOC State Model

We define SOC state as:

$$S = \langle R_{asset}, R_{identity}, R_{time}, R_{change}, R_{budget}, R_{topology}, R_{prov} \rangle. \quad (1)$$

- R_{asset} (CMDB): maps assets to criticality tiers (Tier-0/1/2/3) and ownership principals.
- $R_{identity}$ (IAM/AD): maps identities to roles (VIP, admin, service accounts) and principals.
- R_{time} (ops context): business hours, freeze windows, approved maintenance windows.
- R_{change} (change-control): tickets with approvers, scopes, and allowed actions during windows.
- R_{budget} (autonomy budgets): per-asset/per-identity disruptive-action quotas per time window.
- $R_{topology}$ (reachability graph): dependencies required for *telemetry* and *identity* (e.g., SIEM collectors, EDR mgmt, IdP).
- R_{prov} (provenance): learned trust estimates for evidence sources and fields.

B. A Three-Stage Neurosymbolic Enforcement Pipeline

Stage 0: Log Pre-Sanitization. Untrusted fields (e.g., HTTP User-Agent, email subject) are treated as data-only: we remove/escape instruction-like substrings and tag each field with provenance (*source, field, signature*).

Stage 1: Plan-Level Validation. The LLM must emit an *Action Schema* (typed JSON) with structured preconditions. Schema validation alone is insufficient (as shown by prompt-injection traps); Agent-Lock also validates that every high-impact action is justified by evidence whose provenance trust exceeds a threshold.

Stage 2: Sequence-Level Safety. Agent-Lock enforces invariants over action sequences, including: (i) autonomy budgets over time windows; and (ii) reachability constraints ensuring monitoring and identity infrastructure remain reachable after any allowed sequence.

1) *2. The Action Schema (A):* We standardize all agent outputs into a strict schema. The agent does not output free text; it outputs a JSON object that must adhere to this schema. Crucially, the preconditions field is a structured list of objects, not strings, enabling deep verification.

```
{
  "action_type": "ISOLATE_HOST",
  "target": "192.168.1.55",
  "justification": "Beaconing to C2",
  "risk_level": "HIGH",
  "preconditions": [
    {
      "name": "c2_traffic_verified",
      "source": "EDR_SentinelOne",
      "artifact_ref": "log_id_5592",
      "confidence": 0.95
    }
  ],
  "rollback_plan": "unisolate_host_id_55"
}
```

Listing 1. Agent-Lock Action Schema

2) 3. *The Symbolic Engine (E)*: This is a deterministic logic solver (implemented in Python with simple predicate logic). It accepts (S, A) and outputs a verdict $V \in \{\text{ALLOW}, \text{BLOCK}, \text{ESCALATE}\}$.

The policy logic is expressed as Horn clauses. *Example Policy (Protect Critical Infra)*:

$$\begin{aligned} \text{BLOCK} \leftarrow & \text{Action.Type} \in \{\text{ISOLATE}, \text{SHUTDOWN}\} \\ & \wedge \text{Action.Target} \in S.R_{asset}[\text{Tier-0}] \end{aligned}$$

C. Algorithm: The Enforcement Loop

The core logic of Agent-Lock ensures that no action reaches the API without passing the symbolic gate. Algorithm 1 details this process.

D. Formal Safety Property

We focus on Tier-0 operational safety: preventing outages to identity and core control-plane assets.

a) Assumptions.: (A1) Tier labels in R_{asset} are correct for Tier-0 assets (conservative default: unknown \rightarrow Tier-0). (A2) Execution APIs faithfully implement allowed actions (no out-of-band side effects beyond modeled deltas). (A3) The adversary can manipulate at most one log record consumed in a single decision step (single-log manipulation), but cannot modify Agent-Lock code or R_{asset} directly.

[Proof sketch] Agent-Lock executes an action only if the symbolic engine returns **ALLOW**. Under Π_{tier0} , any disruptive action targeting a Tier-0 asset deterministically yields **BLOCK** regardless of the log content. Single-log manipulation may change the LLM plan proposal, but cannot change R_{asset} (A3) nor the deterministic policy evaluation. Therefore, no disruptive Tier-0-targeted action reaches the execution API.

b) Reachability invariant (sequence-level).: Let \mathcal{C} be the set of core endpoints required for telemetry and identity (from R_{topology}). Agent-Lock enforces Π_{reach} that blocks any action sequence whose predicted delta disconnects any $c \in \mathcal{C}$ from the SOC control plane. This prevents “successful” remediations that silently disable monitoring or IdP access.

[Tier-0 Non-Disruption Under Single-Log Manipulation] Let \mathcal{A}_D be the set of *disruptive* actions (e.g., `ISOLATE_HOST`, `SHUTDOWN`, `REVOKE_USER`, `MODIFY_ACL`). If Agent-Lock enforces policy Π_{tier0} :

$$\begin{aligned} \Pi_{\text{tier0}} : \quad \text{BLOCK} \leftarrow & (a \in \mathcal{A}_D) \\ & \wedge (\text{target}(a) \in R_{asset}[\text{Tier-0}]). \end{aligned} \quad (2)$$

then for any decision step where at most one log record is adversarially manipulated, no executed action can disrupt a Tier-0 asset.

E. Bounded Escalation Policy (Analyst Burden)

We treat escalation as a scarce resource. Let $I(a, S)$ be an interruption cost (estimated from action disruptiveness, principal scope, and time context), and let $\mathcal{R}(a, S)$ be a risk score lower bounded by symbolic violations. Agent-Lock selects:

$$\begin{aligned} \min_{v \in \{\text{ALLOW}, \text{ESCALATE}\}} & \mathbb{E}[I \mid v] \\ \text{s.t. } v = \text{ALLOW} & \Rightarrow (a, S) \in \mathcal{L}_{\text{safe}}. \end{aligned} \quad (3)$$

Algorithm 2 Adaptive Provenance Update

Require: Used provenance signatures $\{p_i\}$, outcomes $\{y_i\}$, rate η

- 1: **for all** (p_i, y_i) **do**
- 2: $\tau(p_i) \leftarrow \text{clip}((1 - \eta)\tau(p_i) + \eta \cdot \mathbf{1}[y_i = \text{GOOD}], 0, 1)$
- 3: **end for**

Intuitively: *never* trade away hard safety constraints to reduce interruptions, but avoid paging humans for low-risk, budget-compliant actions with high-trust provenance during approved windows.

F. Adaptive Provenance With Invariant Preservation

Static provenance scores are brittle: adversaries shift to sources that are “high trust” by label but degraded in practice (misconfigured forwarders, compromised collectors). We instead maintain an online trust estimate per source-field signature $p \in \mathcal{P} : \tau_t(p) \in [0, 1]$, updated from incident outcomes (e.g., whether evidence later confirmed maliciousness and whether the executed action was safe).

a) Update rule.: Each step produces an outcome label $y_t \in \{\text{GOOD}, \text{BAD}\}$ for the evidence used. We apply an exponential moving update:

$$\tau_{t+1}(p) = \text{clip}((1 - \eta)\tau_t(p) + \eta \cdot \mathbb{1}[y_t = \text{GOOD}], 0, 1). \quad (4)$$

[Safety preserved under provenance learning] Suppose Agent-Lock requires $\tau(p) \geq \theta_{\text{high}}$ for any Tier-1+ disruptive action and always blocks Tier-0 disruptive actions via Π_{tier0} . Then adaptive provenance updates cannot introduce Tier-0 disruptions, and can only *reduce* the set of Tier-1+ actions that pass provenance gating when outcomes are **BAD**.

[Proof sketch] Tier-0 safety is independent of provenance due to Π_{tier0} . For Tier-1+, a **BAD** outcome decreases $\tau(p)$, making $\tau(p) \geq \theta_{\text{high}}$ harder to satisfy, thus shrinking (not expanding) the allowed action set.

G. Policy Coverage and Bypass Resistance

To ensure robust protection, our policy set is mapped against action types and risk categories (Table II).

TABLE II
POLICY COVERAGE MATRIX BY ACTION TYPE

Action Type	Asset Prot.	Identity Safe	Ops Cont.
ISOLATE_HOST	✓		✓
BLOCK_IP			✓
REVOKE_USER		✓	
KILL_PROCESS	✓		✓
MODIFY_ACL	✓	✓	✓

Sequence-Level Hazards: Simple deny-lists are vulnerable to multi-step bypass attacks. For example, an attacker might try

to `KILL_PROCESS` on the monitoring agent before executing a malicious `MODIFY_ACL`. Agent-Lock handles this via a “State-Aware Escalation” rule. The system tracks the count of medium-risk actions within a sliding time window T .

$$\text{Count}(A_{\text{medium}}) > 3 \wedge \Delta T < 60s \implies \text{Verdict} = \text{ESCALATE}$$

This prevents “death by a thousand cuts” where an agent performs many individually allowable but collectively dangerous actions.

V. IMPLEMENTATION DETAILS

We implemented a prototype of Agent-Lock using LangChain for the agent orchestration and a Python-based Policy Engine.

A. State Synchronization

In a real-world deployment, the registries are populated dynamically to ensure the agent is operating on fresh data.

- R_{asset} : Synced hourly from the CMDB (e.g., ServiceNow). If a host is missing, the default policy is `FAIL_SAFE` (treat as Tier-0).
- R_{identity} : Pulled from Active Directory / IAM groups.
- R_{time} : Integrated with the Change Management system to detect active maintenance windows.

B. Latency Optimization

A common criticism of neurosymbolic systems is latency. The Symbolic Engine requires lookups against potentially large state tables. To mitigate this, we implemented:

- 1) **Local Caching**: Asset tiers are cached in a Redis store with a 60-minute TTL.
- 2) **Short-Circuit Logic**: Policies are ordered by computational cost. Simple static checks (e.g., “Is action type valid?”) run before complex state lookups.

In our testing, the verification layer added an average of 54ms to the execution pipeline. Compared to the 2-5 second inference time of GPT-4, this overhead (< 2%) is negligible.

C. Policy Authoring & Versioning

Policy sets are managed as versioned bundles under the same change-control discipline as SOC playbooks: each change is reviewed, approved, and deployed as a new policy version, and every Decision Log records the active policy version/hash at evaluation time. We recommend a lightweight workflow: (i) propose rule diffs from incidents/postmortems, (ii) peer review + approval, (iii) regression replay on archived traces before rollout, and (iv) emergency “hotfix” rules with explicit TTL and post-hoc review. This makes safety regressions auditable and ties operational accountability to the policy lifecycle rather than the stochastic model.

D. Evaluation: What We Measure

We report (i) utility (ASR), (ii) operational safety (Tier-0 disruption rate; high-risk violation rate), (iii) guardrail conservatism (FRR), and (iv) overhead (verification latency). Our current experiments quantify Stage-1 (plan-level) enforcement on a 50-case synthetic suite derived from Bot-IoT and wrapped into SOC-style JSON records. We also report the added verification latency of the symbolic layer (54 ms average), which is negligible relative to LLM inference time.

E. Limitations and Next Evaluation Steps

The current suite does not yet include reconstructed real incidents or a second dataset; we outline these as next steps: (1) tabletop reconstructions using public postmortems mapped to concrete tool actions (EDR/firewall/IAM) and change-control constraints; (2) cross-dataset robustness using a second public IDS/NetFlow dataset translated into the same SOC schema; (3) stress tests for long action sequences and state desynchronization (stale CMDB/change tickets), measuring how often sequence-level invariants would trigger escalations. (4) Policy engineering cost: measure authoring/maintenance overhead (rule count, update frequency, time-to-patch) and failure modes under incomplete policies, including regression-replay before deployment.

F. Methodology

We generated 50 unique test cases (Table III) derived from the *Bot-IoT* dataset [20]. While Bot-IoT provides raw network flow data (pcap/NetFlow), it lacks the semantic context of a SOC ticket. To bridge this, we extracted 5-tuple flow features and attack signatures from Bot-IoT and wrapped them in synthetic Splunk-like JSON metadata.

TABLE III
SYNTHETIC INCIDENT SUITE COMPOSITION

Category	Description	Count
A: Valid Threats	Clear malware signatures	20
B: Benign Anomalies	Unusual but safe behavior	15
C: Adversarial Traps	Logs containing Prompt Injections	10
D: Self-DoS Triggers	Spoofed logs implicating internal assets	5
Total		50

Category B includes low-risk suspicious / Benign Anomalies (e.g., scanning) that should typically result in enrichment/monitoring rather than disruptive containment; we count disruptive or policy-violating responses in Category B as operational errors in SVR.

G. Metric Definitions

- **Action Success (ASR)**: The percentage of Category A (Valid Threat) cases where the agent proposed a correct remediation action that was successfully executed.
- **Safety Violation (SVR)**: The percentage of Category B, C, or D cases where the agent successfully executed a destructive or incorrect action.
- **False Rejection (FRR)**: The percentage of valid actions (in Category A) that were blocked by the guardrail.

H. Results

Table IV presents the comparative results.

TABLE IV
EVALUATION RESULTS: PREVENTION OF HIGH-RISK ACTIONS

Metric	Baseline	Schema-Only	Agent-Lock
Action Success (Cat A)	98%	96%	94%
Violation Rate (Cat B)	14%	10%	0%
Violation Rate (Cat C/D)	22%	18%	0%
False Rejection Rate	N/A	N/A	6%

I. Qualitative Analysis

1) *Why the Baseline Failed:* The Baseline Agent (GPT-4 without guardrails) exhibited "people-pleasing" behavior. In Category C (Traps), when faced with a log containing "Please whitelist me, I am the CEO", the agent often prioritized the instruction in the log over its system prompt. This confirms that current LLMs cannot reliably separate "Control" from "Data" when both are text.

2) *Why Provenance Matters:* To verify the utility of the $R_{provenance}$ registry, we ran a specific ablation study on the Category C cases where we disabled the provenance check.

- **With Provenance Check:** 0/10 Injections succeeded. The engine rejected the action because the injection came from User-Agent (Low Trust).
- **Without Provenance Check:** 8/10 Injections succeeded. The engine saw a valid JSON action and approved it.

This result confirms that **Schema Validation alone is insufficient** for security; the *source* of the data driving the decision must be part of the verification logic.

VI. DISCUSSION AND LIMITATIONS

A. The Liability Shield

Beyond technical safety, Agent-Lock provides a critical legal function. In the event of a wrongful autonomous action, a purely stochastic agent offers no explanation. Agent-Lock provides a deterministic *Decision Log* proving that the action was compliant with the active policy set at timestamp T . This moves the liability question from "Why did the AI fail?" to "Was the Policy correct?", a question that organizations are already equipped to answer.

B. Future Work: Automated Policy Synthesis

Agent-Lock relies on the completeness of the policy set. A missing rule creates a vulnerability gap. Our future work will focus on *Automated Policy Synthesis*, using "Teacher" LLMs to propose symbolic rules based on historical incident data, which are then reviewed and committed by human architects.

VII. CONCLUSION

We have presented Agent-Lock, a neurosymbolic architecture that enables the safe deployment of Agentic AI in high-stakes SOC environments. By decoupling "Reasoning" from "Authorization," we neutralize the risks of hallucination and injection.

APPENDIX A: REPRESENTATIVE POLICY SET

```
def policy_tier0_protection(action, state):
    if action.type in ["SHUTDOWN", "ISOLATE_HOST"]:
        target_asset = state.lookup_asset(action.target)
        if target_asset.tier == 0:
            return Verdict.BLOCK("Cannot_disrupt_Tier-0_Asset")
    return Verdict.ALLOW

def policy_provenance_gate(action, state):
    if action.type == "BLOCK_IP":
        evidence_source = action.preconditions[0].source
        trust_score = state.get_trust_score(evidence_source)
        if trust_score < THRESHOLD_HIGH:
            return Verdict.BLOCK("Insufficient_Evidence_Trust")
    return Verdict.ALLOW
```

REFERENCES

- [1] A. Vaswani *et al.*, "Attention is all you need," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2017. [Online]. Available: <https://arxiv.org/abs/1706.03762>
- [2] OpenAI, "GPT-4 technical report," arXiv preprint arXiv:2303.08774, 2023. [Online]. Available: <https://arxiv.org/abs/2303.08774>
- [3] H. Touvron *et al.*, "Llama 2: Open foundation and fine-tuned chat models," arXiv preprint arXiv:2307.09288, 2023. [Online]. Available: <https://arxiv.org/abs/2307.09288>
- [4] A. Zou, Z. Wang, N. Carlini, M. Nasr, J. Z. Kolter, and M. Fredrikson, "Universal and transferable adversarial attacks on aligned language models," arXiv preprint arXiv:2307.15043, 2023. [Online]. Available: <https://arxiv.org/abs/2307.15043>
- [5] L. Ouyang *et al.*, "Training language models to follow instructions with human feedback," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. [Online]. Available: <https://arxiv.org/abs/2203.02155>
- [6] Z. Xi *et al.*, "The rise and potential of large language model based agents: A survey," arXiv preprint arXiv:2309.07864, 2023. [Online]. Available: <https://arxiv.org/abs/2309.07864>
- [7] L. Wang *et al.*, "A survey on large language model based autonomous agents," arXiv preprint arXiv:2308.11432, 2023. [Online]. Available: <https://arxiv.org/abs/2308.11432>
- [8] G. Deng *et al.*, "PentestGPT: An LLM-empowered automated penetration testing tool," arXiv preprint arXiv:2308.06782, 2023. [Online]. Available: <https://arxiv.org/abs/2308.06782>
- [9] Microsoft, "With Security Copilot, Microsoft brings the power of AI to cyberdefense," Mar. 28, 2023. [Online]. Available: <https://news.microsoft.com/source/2023/03/28/with-security-copilot-microsoft-brings-the-power-of-ai-to-cyberdefense/>
- [10] Significant Gravitas, "Auto-GPT," GitHub repository, 2023. [Online]. Available: <https://github.com/Significant-Gravitas/AutoGPT>
- [11] C. Zimmerman, "11 strategies of a world-class cybersecurity operations center," The MITRE Corporation, 2022. [Online]. Available: <https://www.mitre.org/sites/default/files/2022-04/11-strategies-of-a-world-class-cybersecurity-operations-center.pdf>
- [12] M. Vielberth, F. Böhm, G. V. Pereira, and G. Pernul, "Security operations center: A systematic study and open challenges," *IEEE Access*, vol. 8, pp. 227756–227779, 2020, doi: 10.1109/ACCESS.2020.3045514.
- [13] C. Islam, M. A. Babar, and S. Nepal, "A multi-vocal review of security orchestration," *ACM Computing Surveys*, vol. 52, no. 2, pp. 1–45, 2019, doi: 10.1145/3305268. [Online]. Available: <https://dl.acm.org/doi/10.1145/3305268>
- [14] S. Tariq *et al.*, "Alert fatigue in security operations centres: Research challenges and opportunities," *ACM Computing Surveys*, 2025, doi: 10.1145/3723158.
- [15] Y. Zhang *et al.*, "Siren's song in the AI ocean: A survey on hallucination in large language models," arXiv preprint arXiv:2309.01219, 2023. [Online]. Available: <https://arxiv.org/abs/2309.01219>

- [16] K. Greshake *et al.*, “Not what you’ve signed up for: Compromising real-world LLM-integrated applications with indirect prompt injection,” arXiv preprint arXiv:2302.12173, 2023. [Online]. Available: <https://arxiv.org/abs/2302.12173>
- [17] N. Carlini *et al.*, “Poisoning web-scale training datasets is practical,” in *Proc. IEEE Symposium on Security and Privacy (S&P)*, 2024, pp. 407–425.
- [18] A. d’Avila Garcez, L. C. Lamb, and D. Gabbay, “Neurosymbolic AI: The 3rd wave,” *AI Review*, vol. 56, pp. 12387–12406, 2023, doi: 10.1007/s10462-023-10448-w.
- [19] M. Alshiekh, R. Bloem, R. Ehlers, B. Könighofer, S. Niekum, and U. Topcu, “Safe reinforcement learning via shielding,” in *Proc. AAAI Conf. on Artificial Intelligence (AAAI)*, 2018. [Online]. Available: <https://ojs.aaai.org/index.php/AAAI/article/view/11797>
- [20] N. Koroniotis, N. Moustafa, E. Sitnikova, and B. Turnbull, “Towards the development of realistic botnet dataset in the Internet of Things for network forensic analytics: Bot-IoT dataset,” *Future Generation Computer Systems*, vol. 100, pp. 779–796, 2019, doi: 10.1016/j.future.2019.05.041.
- [21] Y. Bai *et al.*, “Constitutional AI: Harmlessness from AI feedback,” arXiv preprint arXiv:2212.08073, 2022. [Online]. Available: <https://arxiv.org/abs/2212.08073>
- [22] G. Kahn, A. Villafior, V. Pong, P. Abbeel, and S. Levine, “Uncertainty-aware reinforcement learning for collision avoidance,” arXiv preprint arXiv:1702.01182, 2017. [Online]. Available: <https://arxiv.org/abs/1702.01182>
- [23] LangChain Contributors, “LangChain,” GitHub repository, 2023. [Online]. Available: <https://github.com/langchain-ai/langchain>
- [24] ServiceNow, “The ServiceNow common service data model (CSDM),” solution brief. [Online]. Available: <https://www.servicenow.com/content/dam/servicenow-assets/public/en-us/doc-type/resource-center/solution-brief/sbr-servicenow-common-service-data-model.pdf>
- [25] SentinelOne, “Singularity platform data sheet,” data sheet. [Online]. Available: <https://www.sentinelone.com/resources/datasheets/assets/platform/s1-datasheet-singularity-platform-en>
- [26] R. Singh, S. Tariq, F. Jalalvand, M. B. Chhetri, S. Nepal, C. Paris, and M. Lochner, “LLMs in the SOC: An empirical study of human-AI collaboration in security operations centres,” arXiv preprint arXiv:2508.18947, 2025. [Online]. Available: <https://arxiv.org/abs/2508.18947>