

ML LABS INTELLIGENCE

Hardening a Production AI Platform's Security Posture

An external audit surfaced 18+ findings on a production AI platform handling sensitive data. ML LABS triaged by severity and remediated in disciplined waves — zero downtime.

CASE STUDY

Author Omar Trejo

Date 2026-03-09

ML LABS

mlabs.com/intel/production-ai-security-hardening

A regulated SaaS platform processing sensitive data across multiple organizations underwent an external security audit. The system had been in production for months — real workloads, real users, real data flowing through AI inference pipelines. The audit produced 18+ findings spanning authentication, input handling, session management, data redaction, and network configuration. None were exotic. Every one was a known vulnerability class that had accumulated under delivery pressure.

ML LABS executed a systematic hardening sprint: triage by severity and blast radius, remediate in priority waves, verify each fix under test. The platform's security posture was transformed without disrupting production operations or requiring downtime.

The Starting Point

The findings clustered in predictable categories, each containing multiple issues that compounded to create a significant attack surface.

- Weak credential storage vulnerable to brute-force attacks
- No rate limiting on authentication endpoints
- Session tokens stored with excessive exposure windows
- Missing security headers on all API responses
- Unsanitized content rendering exposing injection attack surfaces
- XML parsing without injection protection
- Sensitive data redaction that failed silently when the schema expanded

How the Hardening Sprint Worked

ML LABS triaged every finding by blast radius, remediated in priority order, and verified each fix individually. No finding was marked closed until a corresponding test passed — a test designed to fail against the original vulnerability and pass only with the fix applied.

18+ audit findings
across 6 categories



Triage by severity
and blast radius



Disciplined
remediation



Per-fix
verification



Verified and hardened
platform

The ordering mattered. Authentication gates everything else — a system with perfect input validation is still compromised if an attacker can authenticate as a legitimate user. Data handling determines what leaks if other defenses fail. Network configuration locks the perimeter after the interior is sound.

Authentication and Session Hardening

The highest blast-radius findings involved credential storage, session management, and authentication failure handling. ML LABS hardened each surface: credential storage was upgraded to industry-standard adaptive methods, session exposure was reduced to the minimum necessary window, and rate limiting was enforced across all authentication endpoints.

The most critical fix addressed silent authentication failures. The original validation code could fail without raising an error, allowing downstream processing to continue in an unauthenticated context. ML LABS rewrote every validation path to fail explicitly — ambiguity in authentication state was eliminated entirely. Error handling throughout the auth layer was tightened so that security-relevant failures could never be silently swallowed.

The most dangerous authentication code is not code that rejects bad input – it is code that returns nothing and lets the caller decide what nothing means.

Input Handling and Data Redaction

All content rendering surfaces — dashboards, configuration interfaces, and report viewers — were hardened against injection attacks. XML parsing endpoints received injection protection.

The data redaction overhaul was the most consequential change. The existing approach failed silently every time the schema expanded: new fields containing sensitive data would pass through unfiltered. ML LABS inverted the entire redaction model so that nothing passes through unless explicitly marked safe. New fields are blocked by default until reviewed. The failure mode shifted from silent data leakage to overly aggressive filtering — a problem that surfaces immediately rather than compounding invisibly.

Network and Perimeter Controls

Security headers were added across all API responses. CORS policies were tightened to allow only specific production origins. Open redirect vulnerabilities were closed. Per-organization authentication isolation ensured that automated credentials could not be reused across tenants, and cross-tenant access paths were eliminated. Infrastructure configuration was consolidated into a single source of truth with automated drift detection, preventing manual changes from weakening the perimeter after hardening. A web application firewall provided the outer perimeter layer.

When Hardening Is Not Enough

Structured hardening assumes the system is architecturally sound and the findings are implementation-layer issues — wrong hash function, missing header, permissive configuration. When the findings are architectural — no authentication layer at all, no encryption in transit, no tenant data separation — the fix is not a hardening sprint. It is a redesign. Attempting to bolt controls onto a system with fundamental architectural gaps produces a false sense of security while the real exposure remains.

The honest assessment is sometimes that the system needs to be rebuilt with security as a design constraint. The organizations that close their security backlogs fastest recognize this distinction early and concentrate the work in a single operator who holds the full security picture — authentication, data flow, network perimeter, infrastructure configuration — rather than distributing individual findings across a feature team where each engineer sees one piece and none sees the whole.

Results

All 18+ findings were remediated and verified. The redaction model inversion eliminated the class of silent leakage bugs entirely. Per-organization auth isolation was achieved across all tenant boundaries. Zero production downtime occurred during the entire remediation cycle.

- 18+ security findings remediated and individually verified
- Inverted redaction model eliminated silent data leakage
- Per-organization authentication isolation across all tenants
- Zero production downtime during the full remediation cycle
- Automated drift detection prevents configuration regression

First Steps

1. **Audit every surface that handles sensitive data.** Cover authentication, input handling, data redaction, and network configuration. The output should be a prioritized finding inventory classified by blast radius, not a pass/fail report.
2. **Fix authentication first.** Auth weaknesses amplify every other vulnerability class. Harden credential storage, enforce rate limiting on auth endpoints, reduce session exposure, and ensure every validation path fails explicitly rather than silently.
3. **Invert your data redaction model.** Define the set of fields explicitly approved for each output context and block everything else by default. This single change eliminates the most dangerous category of silent failures.

Practical Solution Pattern

Execute security hardening as a structured sprint with triage-by-severity, disciplined remediation ordering, and per-fix verification. Prioritize authentication because it gates everything else, then close data redaction gaps by inverting the default from "allow unless blocked" to "block unless allowed," and lock the network perimeter with explicit controls and automated drift detection. Every fix ships with a test that would have caught the original vulnerability.

This works because production AI security failures are overwhelmingly known vulnerability classes with known fixes — the gap is execution discipline, not knowledge. Organizations that concentrate remediation authority in a single experienced operator who holds the full security picture — authentication, data flow, network configuration, infrastructure — close their backlogs in tight cycles rather than watching findings accumulate across distributed teams. Deep expertise paired with AI-augmented execution now matches or exceeds what large security

teams produce, with less coordination overhead and faster iteration through the audit-triage-remediate-verify loop. If your platform handles sensitive data and has not undergone a structured security pass, a **Technical AI Assessment** can surface the findings before an external audit does.

References

1. HackerOne. **AI Security Findings Outpace Cybersecurity Remediation in 2025**. *HackerOne*, 2025.
2. OWASP. **Top 10 for LLM Applications 2025**. *OWASP Foundation*, 2025.
3. NIST. **Cybersecurity Framework Profile for Artificial Intelligence (NIST IR 8596)**. *NIST*, 2025.
4. OWASP. **Password Storage Cheat Sheet**. *OWASP Cheat Sheet Series*, 2024.
5. OWASP. **Cross Site Scripting Prevention Cheat Sheet**. *OWASP Cheat Sheet Series*, 2024.
6. NIST. **Healthcare Security Rule Guidance (SP 800-66)**. *NIST*, 2024.



ML LABS

Custom AI Systems for High-Value Workflows

mllabs.com