

AMTSO Standard Compliance Statement

This Test has been designed to comply with the Anti-Malware Testing Standards Organization, Inc. ("AMTSO") Testing Protocol Standard for the Testing of Anti-Malware Solutions, Version [1.3] (the "Standard"). This Test Plan has been prepared using the AMTSO Test Plan Template and Usage Directions, Version [2.4]. SecureQLab is solely responsible for the content of this Test Plan.



SecureIQLab

Methodology

AI Security CyberRisk Validation Methodology v1.0

Version:	1.0
Last Revision:	4 th March, 2026
Language:	English

Table of Contents

1	INTRODUCTION	2
1.1	DEFINITION	2
1.2	AI SECURITY VS WEB APPLICATION FIREWALL (WAF)	2
1.3	AI THREAT LANDSCAPE AND EMERGING RISK	3
1.4	ENTERPRISE ADOPTION DRIVERS	3
1.5	RATIONALE FOR INDEPENDENT CYBER RISK VALIDATION	4
2	SCOPE AND OBJECTIVES	4
2.1	SCOPE	4
2.2	OBJECTIVES	5
3	CYBER RISK VALIDATION CRITERIA	5
3.1	SECURITY EFFICACY	5
3.1.1	Input Security	5
3.1.2	Output Security	6
3.1.3	Retrieval Firewall	8
3.1.4	False Positive Testing	8
3.2	OPERATIONAL EFFICIENCY	9
3.2.1	Deployment & Onboarding	9
3.2.2	Policy Management & Administration	9
3.2.3	Integration with Enterprise Ecosystem	10
3.2.4	Incident Response and Visibility	10
3.2.5	Insight for Threat Hunting and Forensics	10
3.2.6	Security Administration	10
4	VALIDATION METHODOLOGY	11
4.1	TEST ENVIRONMENT PREPARATION	11
4.2	BASELINE CONFIGURATION	12
4.3	INITIAL SMOKE TESTING	12
4.4	ADVERSARIAL SECURITY TESTING	12
4.5	OPERATIONAL EFFICIENCY VALIDATION	13
4.6	SCORING AND SCORECARD COMPILATION	13
4.7	VENDOR REVIEW AND DISPUTE RESOLUTION	13
4.8	FINAL REPORT PUBLICATION	14
5	AI SECURITY VALIDATION ARCHITECTURE AND SETUP OVERVIEW	14
6	SCORING MECHANISM FOR CYBERRISK VALIDATION	15
6.1	SECURITY EFFICACY SCORING CRITERIA	15
6.2	OPERATIONAL EFFICIENCY SCORING CRITERIA	15
7	GENERAL EVALUATION APPROACH	18
7.1	AI SECURITY VENDOR PARTICIPATION SELECTION CRITERIA	18
7.2	LIST OF CONSIDERED AI SECURITY VENDORS	18
7.3	VALIDATION TIMELINE	19
7.4	RISK AND RISK MANAGEMENT	20
7.5	GEO LIMITATION	20
7.6	DISTRIBUTION OF TEST DATA	20
7.7	FUNDING AGREEMENT	20
7.8	DISPUTE PROCESS	20
8	OPT-OUT POLICY	21
9	ATTESTATIONS	22
10	DOCUMENT AND REVISION	22
11	COPYRIGHT AND DISCLAIMER	23
12	APPENDICES	24
12.1	VALIDATION FRAMEWORK MAPPING TABLE	24

1 INTRODUCTION

AI is now embedded in enterprise systems that handle customer interactions, software development, and internal data processing. Adoption has rapidly moved from pilot programs to operational use. A 2024 survey by [McKinsey](#) shows that most organizations use AI in at least one business function, with AI based systems increasingly integrated into core workflows and connected to sensitive data.

Business outcomes, however, trail adoption. Research associated with [MIT](#) indicates that only a small share of initiatives deliver measurable impact at scale. Many AI deployments fail to progress into stable production due to unresolved risks related to data exposure, unpredictable model behavior, regulatory responsibility, and unclear control ownership.

Taken together, these findings illustrate a clear reality: while AI adoption is advancing faster than any previous wave of enterprise technology, security assurance and governance maturity are not keeping pace. Organizations integrating AI into critical or regulated workflows require specialized security controls, continuous monitoring, and independent risk validation to ensure that these systems are trustworthy, auditable, and safe for production use.

1.1 DEFINITION

SecureIQLab defines AI security as a runtime enforcement and governance layer that mediates all interactions between users, AI models, external data sources, and connected tools. It inspects and controls prompts, system instructions, retrieved context, model outputs, and agent actions to detect and block prompt injection, data exfiltration, policy violations, tool abuse, and unsafe autonomous behavior.

Unlike traditional security controls that operate at the network or application layer, AI security operates at the semantic and decision layer, enforcing intent-aware policies, contextual access controls, and action-level authorization.

Its objective is to reduce model-driven risk by ensuring that AI systems remain aligned with organizational policies, least-privilege principles, and human governance throughout the full execution lifecycle.

1.2 AI SECURITY VS WEB APPLICATION FIREWALL (WAF)

Dimension	AI Security	Web Application Firewall
Primary focus	Governs AI model behavior, prompt handling, retrieval context, outputs, and agent actions	Filters and protects web applications from malicious HTTP traffic
Protection Layer	Semantic and decision layer	Application Layer (Layer 7)
Traffic Type	Natural language, embeddings, system prompts, tool calls	HTTP requests, headers, parameters, payloads
Threat Detection Model	Intent-aware, contextual, probabilistic analysis	Deterministic rules, signatures, and pattern matching
Threat Surface	Prompts, RAG pipelines, memory, system instructions, function calls	APIs, forms, URLs, cookies, query strings
Context Awareness	High (conversation state, intent hierarchy, retrieval context)	Limited (request/session-level awareness)

OWASP coverage	OWASP Top 10 for LLM Applications (2025)	OWASP Top Ten Web Application Security Risks (2025)
----------------	--	---

1.3 AI THREAT LANDSCAPE AND EMERGING RISK

The OWASP Top 10 for LLM Applications (2025) identifies security risks that arise from the probabilistic reasoning, natural-language interfaces, retrieval pipelines, and autonomous execution capabilities of modern AI systems. AI-native threats manifest across the AI execution lifecycle and can be categorized as follows:

- Prompt injection (direct, indirect, multimodal)
- Sensitive information disclosure
- Supply-chain compromise of models and adapters
- Data and model poisoning
- Improper output handling
- Excessive agency in agentic systems
- System prompt leakage
- Vector and embedding weaknesses (RAG abuse)
- Misinformation and hallucination risks
- Unbounded consumption and cost-exhaustion attacks

These risks introduce confidentiality, integrity, availability, financial, and reputational impacts unique to AI systems.

1.4 ENTERPRISE ADOPTION DRIVERS

Enterprises are rapidly integrating AI into core business systems, transitioning from experimentation to production deployment. However, adoption is no longer driven solely by innovation and productivity gains; it is increasingly shaped by the need to manage AI-specific risk at scale.

The following factors are driving enterprise adoption of AI technologies across core business systems and operational workflows.

Productivity Automation

AI technologies are increasingly embedded in customer support, software development, research, and knowledge management workflows. Organizations leverage AI to automate repetitive tasks, accelerate decision-making, and augment human expertise, often achieving significant productivity gains across engineering, operations, and business functions.

Cost Reduction and Operational Efficiency

By automating high-volume, language-intensive processes such as ticket handling, documentation, code generation, and data analysis, enterprises aim to reduce operational overhead, shorten cycle times, and optimize workforce allocation. In many cases, AI applications are positioned as force multipliers rather than direct replacements for human workers.

Competitive Differentiation Through AI-Driven Services

AI enables organizations to deliver differentiated, AI-native products and services, such as intelligent assistants, personalized experiences, and conversational interfaces, that enhance customer engagement and create new revenue opportunities. As AI capabilities become table stakes across multiple industries, the speed of adoption directly impacts competitive positioning.

Agent-Based Orchestration of Workflows

The emergence of agentic architectures, in which AI systems autonomously invoke tools, APIs, and workflows, has

expanded the role of AI systems from passive responders to active system operators. Enterprises increasingly rely on AI agents to coordinate multi-step processes across IT, DevOps, security, finance, and other business systems.

Retrieval-Augmented Generation (RAG) Over Proprietary Data

To unlock value from internal knowledge bases, enterprises deploy RAG architectures that combine LLMs with proprietary documents, databases, and embeddings. This approach enables organizations to apply generative AI to sensitive internal data without fully retraining models, significantly lowering barriers to adoption.

1.5 RATIONALE FOR INDEPENDENT CYBER RISK VALIDATION

AI risks are probabilistic, contextual, and behavior-driven, and they cannot be reliably evaluated through feature checklists, marketing claims, or limited proof-of-concept deployments. Capabilities such as prompt-injection prevention, retrieval security, agent control, and misuse detection must be validated against adversarial techniques, chained attack scenarios, and production-scale workloads.

Independent cyber risk validation provides enterprises with:

- **Vendor-neutral benchmarking** of AI security solutions, assessing security efficacy, operational effectiveness, and enterprise readiness across the [OWASP Top 10 for LLM Applications](#)
- **Reproducible and transparent test methodologies** that simulate real-world adversary behavior, misuse patterns, and failure modes unique to LLM systems.
- **Evidence-based insights** that support procurement decisions, deployment strategies, and risk acceptance discussions with executive leadership, regulators, and boards.

SecureIQLab's AI Security Cyber Risk Validation Methodology ensures that AI security solutions are evaluated not only for their stated features, but also for their ability to measurably reduce cyber risk in production environments. By testing against realistic attack paths, operational constraints, and enterprise integration requirements, the methodology enables organizations to adopt AI security controls with confidence, clarity, and defensible assurance.

2 SCOPE AND OBJECTIVES

2.1 SCOPE

This methodology is designed to evaluate AI security solution capabilities, defined as security controls that operate at or inline with AI interactions to prevent, detect, and respond to AI-specific threats as outlined in the [OWASP Top 10](#) for Large Language Model Applications.

In-Scope Vendors

The validation includes vendors that meet at least one of the following criteria:

- Pure-Play LLM Firewall
- Broader AI Security Solutions
- API Security and Edge Platforms Offering LLM Protection

Out-of-Scope Technologies

The following are explicitly out of scope for this methodology and are not evaluated or scored:

- AI Security Posture Management (AI-SPM) platforms without runtime enforcement
- Data Security Posture Management (DSPM) tools
- CASB, SSE, or browser isolation platforms without in-line LLM controls
- Model training pipelines, fine-tuning workflows, or MLOps platforms
- General-purpose API gateways or WAFs lacking LLM-specific logic
- Standalone Agentic AI Firewalls and MCP Firewalls

Capabilities in these categories may be referenced for context but do not influence validation scores.

2.2 OBJECTIVES

The primary objective of this validation is to provide independent, evidence-based assurance of an AI Security solution's ability to reduce real-world cyber risk while remaining operationally viable in enterprise environments.

A. Quantify Security Efficacy

This objective measures how effectively an AI Security solution detects, prevents, and mitigates AI-specific threats under realistic adversarial conditions.

Key validation goals include:

- **Measure detection and prevention effectiveness** against the OWASP Top 10 for LLM Applications (2025), covering prompt injection, sensitive data disclosure, retrieval abuse, excessive agency, system prompt leakage, misinformation, and unbounded consumption.
- **Evaluate resilience to adversarial and chained attack techniques**, including indirect prompt injection, obfuscated prompts, multi-turn manipulation and retrieval poisoning.
- **Validate decision accuracy**, including analysis of false-positive and false-negative rates, to ensure security enforcement does not disrupt legitimate enterprise workflows or create excessive operational noise.

B. Validate Operational Efficiency & Enterprise Readiness

This objective assesses whether an AI Security solution can be deployed, managed, and operated at enterprise scale without introducing undue complexity or friction.

Key validation goals include:

- **Assess deployment complexity and time-to-value**, including installation models (inline proxy, SDK, gateway), configuration effort, and alignment with enterprise AI architectures.
- **Evaluate manageability and scalability**, including policy administration, multi-tenant support, performance impact, and the ability to scale with production workloads.
- **Measure security operations enablement**, including visibility into LLM interactions, centralized logging, alerting, SOC integration (SIEM/SOAR), and incident response readiness.

Validation Outcome

Together, these objectives ensure that participating AI Security solutions are evaluated not only on their ability to block attacks, but on their practical suitability for sustained enterprise deployment providing CISOs and security leaders with defensible data to support procurement, risk acceptance, and governance decisions.

3 CYBER RISK VALIDATION CRITERIA

The Cyber Risk Validation Criteria define the dimensions against which AI security solutions are independently assessed. These criteria are designed to measure both technical security efficacy and operational readiness under realistic enterprise and adversarial conditions.

3.1 SECURITY EFFICACY

Security efficacy measures an AI Security solution's ability to detect, prevent, and mitigate LLM-specific threats in real time, aligned with the OWASP Top 10 for LLM Applications (2025). Security efficacy is validated through a three-layer enforcement model consisting of Input Security, Retrieval Firewall, and Output Security controls.

3.1.1 Input Security

This security includes all defenses that inspect and secure content flowing from the user to the GenAI Application. The goal is to prevent the user from manipulating, compromising, or abusing the system.

- **Prompt Injection & Jailbreak Detection:** The most fundamental capability. This involves detecting and blocking deceptive user inputs (prompts) designed to bypass safety filters, ignore previous instructions, or cause the model to perform unintended actions.
- **Data and Model Poisoning:** This involves detecting and preventing malicious attempts to corrupt model behavior or downstream knowledge sources through crafted inputs.
- **Sensitive data ingestion,** such as PII, credentials, secrets, or proprietary information submitted via prompts.
- **Obfuscated, multilingual, and multimodal inputs,** including encoding tricks, translation-based bypasses, and embedded instructions.

The following are the test validation scenarios for input security:

Validation Scenario 1: Direct Prompt Injection

- **Objective:** To validate the AI Security's ability to detect and prevent direct prompt injection attempts that attempt to override system instructions or bypass established safety controls.

Validation Scenario 2: Indirect Prompt Injection

- **Objective:** To validate the AI Security's ability to detect and prevent indirect prompt injection attacks delivered through external or embedded content.

Validation Scenario 3: Multilingual / Obfuscated Attack (Direct Prompt Injection)

- **Objective:** To validate the AI Security's ability to detect and prevent direct prompt injection attempts delivered through multilingual phrasing, encoding techniques, character substitution, or other obfuscation methods designed to evade detection controls.

Validation Scenario 4: Vector Database Toxic Content Injection

- **Objective:** To validate whether the AI Security solution prevents harmful content from influencing retrieval-based responses.

Validation Scenario 5: PII Disclosure via Direct User Input

- **Objective:** To validate whether sensitive data ingestion controls detect and block user-submitted PII from being processed, stored, or used in responses.

Validation Scenario 6: Payment Card Data Ingestion (PCI)

- **Objective:** To validate whether sensitive data ingestion controls detect and block payment card information to enforce PCI-related restrictions.

Validation Scenario 7: Oversized Prompt Blocking

- **Objective:** To validate the AI Security solution ability to detect and block excessively large prompts designed to exhaust compute resources.

Validation Scenario 8: High-Volume Request Flood

- **Objective:** To validate whether the AI Security solution enforces request rate limits and per-user quotas.

3.1.2 Output Security

This category includes all defenses that inspect and secure content flowing from the GenAI model to the user. The goal is to prevent the model from exposing sensitive information or generating harmful content.

- **Data Loss Prevention (DLP):** This is the most critical output-side control. It involves real-time detection and blocking of sensitive data exfiltration. This includes Personally Identifiable Information (PII) such as Social Security Numbers, phone numbers, and email addresses; financial data like credit card (Luhn algorithm validated) and IBAN numbers; and user credentials.
- **Intellectual Property (IP) & Sensitive Data Redaction:** A more advanced form of DLP, this capability focuses on protecting the organization's proprietary data, such as internal source code, financial projections, or strategic product roadmaps, from being "leaked" in a model's response.

- **Toxic & Harmful Content Filtering:** This involves filtering model-generated responses for hate speech, violence, sexual content, self-harm, and other forms of inappropriate or off-brand content. Advanced forms of this capability also include detecting and flagging non-grounded "hallucinations" or copyrighted material.

The following are the test case scenarios for Output Security:

Validation Scenario 9: Direct Sensitive Data Extraction

- **Objective:** To validate whether the AI Security solution prevents explicit attempts to extract sensitive information such as PII, credentials, financial records, or confidential business data.

Validation Scenario 10: Cross-Session Data Leakage

- **Objective:** To validate whether the AI Security solution enforces session isolation and prevents disclosure of another user's data.

Validation Scenario 11: Prompt Injection for Sensitive Data Disclosure

- **Objective:** To validate whether the AI Security solution prevents sensitive data disclosure triggered by embedded prompt injection.

Validation Scenario 12: Command Injection via LLM Output

- **Objective:** To validate whether the AI Security solution prevents execution of malicious shell commands generated by the LLM.

Validation Scenario 13: Malicious Content in Email Templates

- **Objective:** To validate sanitization of LLM-generated email content.

Validation Scenario 14: Cross-Site Scripting (XSS) via LLM Output

- **Objective:** To validate whether the application properly encodes LLM-generated HTML or JavaScript content.

Validation Scenario 15: SQL Injection via LLM-Generated Query

- **Objective:** To validate whether LLM-generated SQL queries are safely parameterized before execution.
- **Context:** An LLM is used to generate SQL queries from natural language. If queries are executed directly, destructive operations may occur.

Validation Scenario 16: Toxic and Harmful Content

- **Objective:** To validate the AI Security solution ability to monitor and block the generation of toxic, harmful, or policy-violating content.

Validation Scenario 17: Excessive Agency

- **Objective:** To verify that the AI security solution can detect and prevent unsafe or unintended plugin executions caused by manipulated LLM outputs. The objective is to assess capability-level controls addressing excessive agency risks, not to evaluate full agent orchestration infrastructure or dedicated runtime firewall architectures

Validation Scenario 18: Direct System Prompt Extraction Attempt

- **Objective:** To validate whether the AI Security solution prevents direct attempts to extract system prompt content.

Validation Scenario 19: System Prompt Reconstruction via Systematic Querying

- **Objective:** To validate whether the AI Security solution detects and prevents systematic probing attempts intended to reconstruct or infer system prompt content.

Validation Scenario 20: Guardrail Bypass via Extracted System Prompt

- **Objective:** To validate whether the AI Security solution controls and prevents attackers from extracting system prompt guardrails and using that knowledge to bypass restrictions, potentially leading to remote code execution or other downstream exploitation.

Validation Scenario 21: Fabricated Citation & Reference Detection

- **Objective:** To validate the AI Security solution ability to detect hallucinated academic citations, fabricated references, or unverifiable sources generated by the LLM.

3.1.3 Retrieval Firewall

This category evaluates AI Security controls that govern context retrieval mechanisms, including Retrieval-Augmented Generation (RAG), memory layers, and vector-based knowledge access. The objective is to ensure that external context supplied to the LLM is authorized, trustworthy, relevant, and isolated, preventing retrieval-driven compromise of confidentiality, integrity, or model behavior.

Key Retrieval Firewall Capabilities Evaluated:

- **Vector Isolation and Tenant Separation:** Validation focuses on the firewall's ability to enforce strict logical isolation between users, applications, and tenants within shared vector stores.
- **Poisoned Document Detection:** This capability evaluates the firewall's ability to identify and mitigate malicious or misleading content introduced into retrieval sources.

The following are the test case scenarios for retrieval firewall:

Validation Scenario 22: Vector and Embedding Weakness

- **Objective:** To assess the security of the vector database and embedding processes used in RAG systems.

Validation Scenario 23: Poisoned Document with Hidden Prompt Injection (Indirect Injection in RAG)

- **Objective:** To validate the AI Security's ability to detect and prevent indirect prompt injection originating from retrieved RAG content, including hidden or obfuscated instructions embedded inside documents, web pages, PDFs, or metadata.

Validation Scenario 24: Poisoned RAG Upload With "Misinformation Injection" (Decision Manipulation)

- **Objective:** To validate the AI Security's ability to detect, prevent, and mitigate misinformation injected into Retrieval-Augmented Generation (RAG) sources that are intended to manipulate business decisions, policies, or operational outcomes.

3.1.4 False Positive Testing

False Positive Validation evaluates whether the AI Security Control Layer blocks, rewrites, redacts, or escalates legitimate enterprise use-cases that are superficially similar to attack patterns. The goal is to measure precision, business disruption risk, and policy usability under realistic workloads.

The following are the test case scenarios for false positive testing:

Validation Scenario 25: Security Awareness Content (Benign Prompt Injection Discussion)

- **Objective:** To validate that the AI Security Control Layer does not block legitimate security training content that references prompt injection or jailbreak techniques in a non-exploitative context.

Validation Scenario 26: Quoted Malicious Text for Analysis (Treat as Data, Not Instruction)

- **Objective:** To validate that the AI Security Control Layer does not block or misclassify malicious-looking text when it is clearly provided for analysis, classification, or defensive review rather than execution.

Validation Scenario 27: Benign Business Identifiers Mistaken as Credentials

- **Objective:** To validate that the AI Security Control Layer does not incorrectly classify common business identifiers (e.g., ticket numbers, order IDs, invoice numbers) as sensitive credentials or secrets requiring blocking or redaction.

Validation Scenario 28: Benign DevOps Command Explanation (Non-Executable Context)

- **Objective:** To validate that the AI Security Control Layer does not misclassify legitimate operational troubleshooting requests as command injection attempts when the user is requesting explanation or education rather than execution.

Validation Scenario 29: Multilingual Benign Prompt with Mixed Scripts

- **Objective:** To validate that the AI Security Control Layer does not elevate risk scores, block, or degrade service for legitimate multilingual prompts, including those written entirely in non-English languages or containing mixed scripts.

Validation Scenario 30: PCI-Like Numeric String in Non-Payment Context

- **Objective:** To validate that the AI Security Control Layer does not automatically classify any 16-digit numeric sequence as payment card data (PCI) without contextual validation and checksum verification.

Validation Scenario 31: Legitimate Policy Content Discussing Sensitive or Restricted Topics

- **Objective:** To validate that the AI Security Control Layer does not block or incorrectly flag legitimate compliance, HR, or legal content solely because it contains terms commonly associated with restricted or sensitive categories (e.g., harassment, violence, abuse, discrimination).

Validation Scenario 32: Large Document Summarization (Legitimate High-Token Input)

- **Objective:** To validate that the AI Security Control Layer does not misclassify legitimate large inputs (e.g., policies, contracts, technical manuals, audit reports) as resource exhaustion or denial-of-service attempts solely due to token length.

Note: For details regarding Test Scenarios, please contact SecureIQLab.

3.2 OPERATIONAL EFFICIENCY

Operational efficiency evaluates whether an AI Security solution can be deployed, managed, and operated reliably at enterprise scale without introducing excessive complexity, latency, or operational burden. While security efficacy measures what the firewall can stop, operational efficiency determines whether it can be sustained in production environments.

3.2.1 Deployment & Onboarding

This criterion assesses the practicality, flexibility, and time-to-value of deploying the AI Security solution into real-world enterprise architectures.

Validation focuses on:

- **Deployment models supported**, including inline proxy, API gateway, SDK-based instrumentation, sidecar, or service mesh integration
- **Compatibility with enterprise environments**, including cloud-native (AWS, Azure, GCP), hybrid, and on-prem deployments
- **Support for common LLM architectures**, such as chat-based applications, RAG pipelines, and agent-enabled workflows
- **Initial onboarding effort**, including configuration complexity, required code changes, dependency management, and setup documentation
- **Time-to-value**, measured as the time required to reach a functional, policy-enforced state using recommended default configurations

3.2.2 Policy Management & Administration

This criterion evaluates the effectiveness and usability of the AI Security's policy framework.

Validation focuses on:

- **Policy authoring models**, including rule-based, ML-assisted, or hybrid approaches
- **Policy granularity**, enabling differentiated controls based on user role, application, intent, data type, retrieval source, or action type
- **Policy explainability**, ensuring administrators can understand why a prompt, retrieval, or output was allowed, blocked, or modified
- **Change management**, including policy versioning, staged rollout, rollback, and audit history

- **Separation of duties**, allowing different roles (security, platform, application owners) to manage policies without excessive privilege

3.2.3 Integration with Enterprise Ecosystem

This criterion assesses how well the AI Security solution integrates with existing enterprise security, IT, and operations tooling.

Validation includes:

- **Security platform integration**, such as SIEM, SOAR, and XDR solutions for alert ingestion and response automation
- **Identity and access integration**, including IAM, SSO, RBAC, and attribute-based access control for identity-aware enforcement
- **Data protection integration**, including interoperability with DLP or data classification systems
- **API and telemetry export**, supporting real-time streaming, batch export, or API-based access to events and metrics

3.2.4 Incident Response and Visibility

This criterion measures how effectively the AI Security solution enables rapid detection and response to LLM-related security incidents.

Validation focuses on:

- Real-time alerting, with actionable alerts tied to specific LLM threats (e.g., prompt injection, retrieval poisoning, data leakage)
- Incident context reconstruction, including visibility into prompts, retrieved context, model outputs, tool calls, and enforcement actions
- Evidence preservation, ensuring that logs, artifacts, and metadata are retained in a forensically sound manner
- Correlation capabilities, linking LLM events with broader security telemetry (user identity, device, network, application)

3.2.5 Insight for Threat Hunting and Forensics

This criterion evaluates advanced analytical capabilities that support proactive defense and post-incident analysis.

Validation includes:

- Prompt and response traceability, across sessions, users, and applications
- Conversation-level timelines, enabling reconstruction of multi-step or slow-burn attacks (e.g., gradual extraction or poisoning)
- Attack pattern analytics, identifying recurring techniques such as repeated probing, semantic enumeration, or policy evasion attempts
- Search and filtering capabilities, allowing analysts to query historical LLM interactions by risk type, user, model, or data category

3.2.6 Security Administration

This criterion evaluates the administrative controls required for governance, compliance, and enterprise-scale operations.

Validation focuses on:

- Role-based access control (RBAC) for administrators, analysts, auditors, and operators
- Audit logging, capturing administrative actions, policy changes, and enforcement decisions

- Compliance reporting, including exportable evidence aligned with regulatory and governance requirements (e.g., OWASP LLM Top 10, NIST AI RMF)
- Operational resilience, including backup, configuration recovery, and platform availability controls

4 VALIDATION METHODOLOGY

The methodology is designed to ensure fairness, reproducibility, transparency, and real-world relevance, while minimizing vendor bias and test variability.

The validation process evaluates both security efficacy and operational efficiency using a combination of automated testing, human-led adversarial simulation, and enterprise workflow validation.

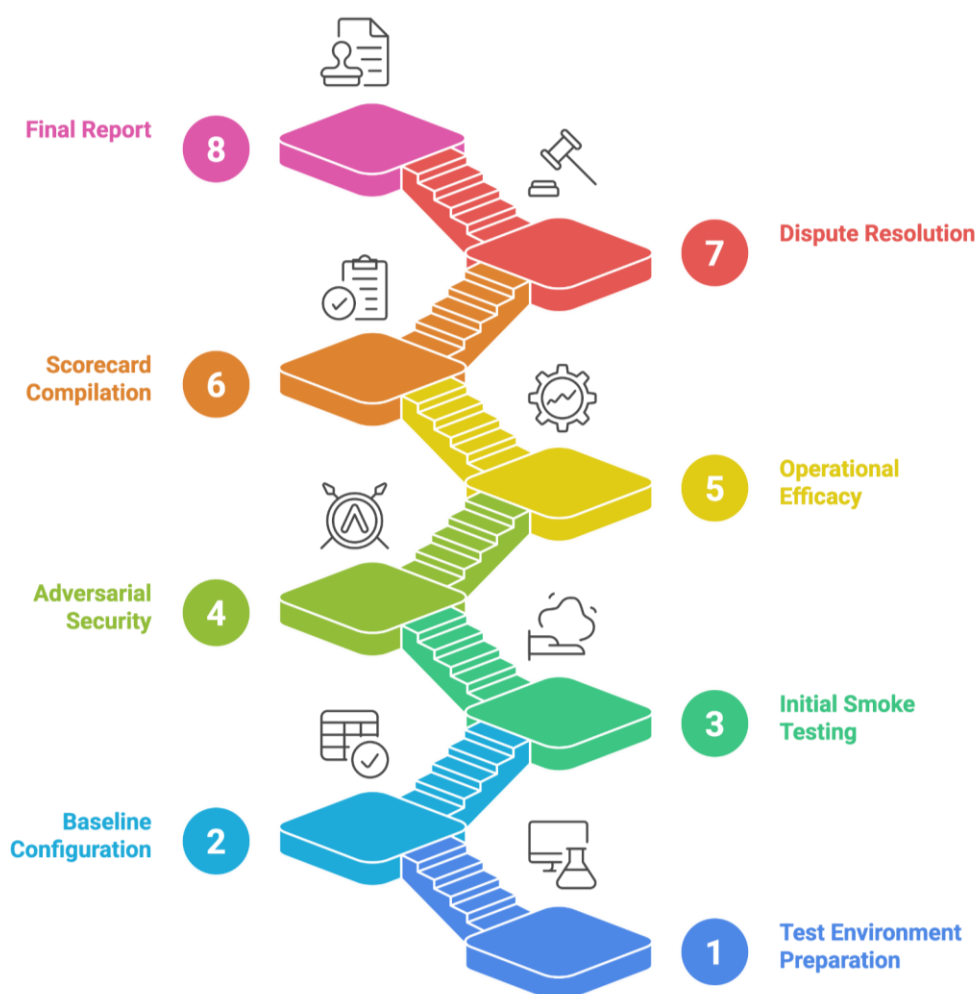


Figure 1. Test validation process

4.1 TEST ENVIRONMENT PREPARATION

SecureQLab establishes a controlled, production-representative test environment to ensure that validation results reflect real-world enterprise deployments.

Environment preparation includes:

- Deployment of one or more LLM-enabled applications.
- Configuration of representative enterprise data sets:
 - Public and internal documents
 - Sensitive and synthetic PII data
 - Canary values for leakage detection
- Integration with identity providers, logging systems, and monitoring tools to simulate enterprise security operations.
- Isolation of each participating vendor's test environment to prevent cross-contamination or shared learning effects

All test data used is synthetic, anonymized, or purpose-built, ensuring no real customer data is exposed.

4.2 BASELINE CONFIGURATION

Each AI Security solution is initially deployed using vendor-recommended default configurations to reflect how enterprises commonly deploy security controls.

Baseline configuration principles include:

- Enablement of all generally available AI Security protections relevant to the scoped domains
- No custom tuning, signatures, or rules beyond documented defaults
- Policy configuration aligned with publicly available vendor guidance
- Validation of successful deployment via smoke testing prior to adversarial evaluation

Any deviations from baseline configuration are documented and disclosed.

4.3 INITIAL SMOKE TESTING

Before executing full adversarial testing, SecureIQLab performs smoke testing to validate functional readiness.

Smoke testing validates:

- Proper interception of LLM inputs, retrieval flows, outputs, and agent actions
- Correct policy enforcement for basic high-risk scenarios (e.g., simple prompt injection, obvious PII leakage)
- Logging, alert generation, and telemetry export functionality
- Stability and performance under nominal load

4.4 ADVERSARIAL SECURITY TESTING

SecureIQLab conducts structured adversarial testing aligned to the OWASP Top 10 for LLM Applications (2025).

Testing characteristics include:

- Execution of predefined test cases across:
 - Input security
 - Output security
 - Retrieval firewall
 - Abuse, misuse, and cost controls
- Use of:
 - Direct and indirect prompt injection
 - Obfuscation, multilingual, and chained attack techniques
 - Retrieval poisoning and semantic leakage attacks
 - Multi-turn and slow-burn extraction attempts
- Combination of:
 - Automated attack scripts for scale and consistency

- Manual red-team testing for creativity and evasive techniques

Each test case is executed multiple iterations to ensure consistency and accuracy.

4.5 OPERATIONAL EFFICIENCY VALIDATION

In parallel with security testing, SecureIQLab evaluates operational characteristics.

Validation includes:

- Measurement of deployment complexity and configuration effort
- Policy authoring, versioning, and rollback workflows
- SOC integration, alert quality, and investigation workflows
- Performance impact assessment, including:
 - Request latency
 - Token throughput
 - Resource utilization overhead

Operational testing ensures security efficacy does not come at the cost of enterprise usability.

4.6 SCORING AND SCORECARD COMPILATION

Results from all test cases are aggregated into a standardized scorecard.

Scoring principles include:

- Separate scoring for:
 - Security efficacy
 - Operational efficiency
- Weighted scoring based on risk impact and enterprise relevance
- Transparent documentation of:
 - Passed tests
 - Partially successful tests
 - Failed tests
- Exclusion of out-of-scope capabilities from scoring

Scoring and Disclosure Treatment

- Vendors offering broader security platforms are scored only on their AI Security components
- Opt-in domains must be declared prior to testing
- Non-participation in specific domains is transparently disclosed in final reports
- No vendor is advantaged or penalized for features outside the defined AI Security scope

4.7 VENDOR REVIEW AND DISPUTE RESOLUTION

Participating vendors are provided with preliminary results for review.

The dispute process includes:

- Controlled review period for vendors to:
 - Validate findings
 - Identify potential misconfigurations
 - Submit evidence-based disputes
- Limited retesting for validated configuration errors
- Equal application of any scoring adjustments across all vendors

SecureIQLab retains final authority over scoring decisions.

4.8 FINAL REPORT PUBLICATION

Upon completion of validation:

- Individual vendor reports are produced
- A comparative report is published
- Methodology, scoring criteria, and limitations are fully disclosed

If previously unknown security vulnerabilities are identified during testing:

- Findings are disclosed privately to the affected vendor
- Vendors are provided a remediation window prior to public disclosure
- Public reporting follows responsible disclosure best practices

This ensures research integrity while minimizing unnecessary risk.

5 AI SECURITY VALIDATION ARCHITECTURE AND SETUP OVERVIEW

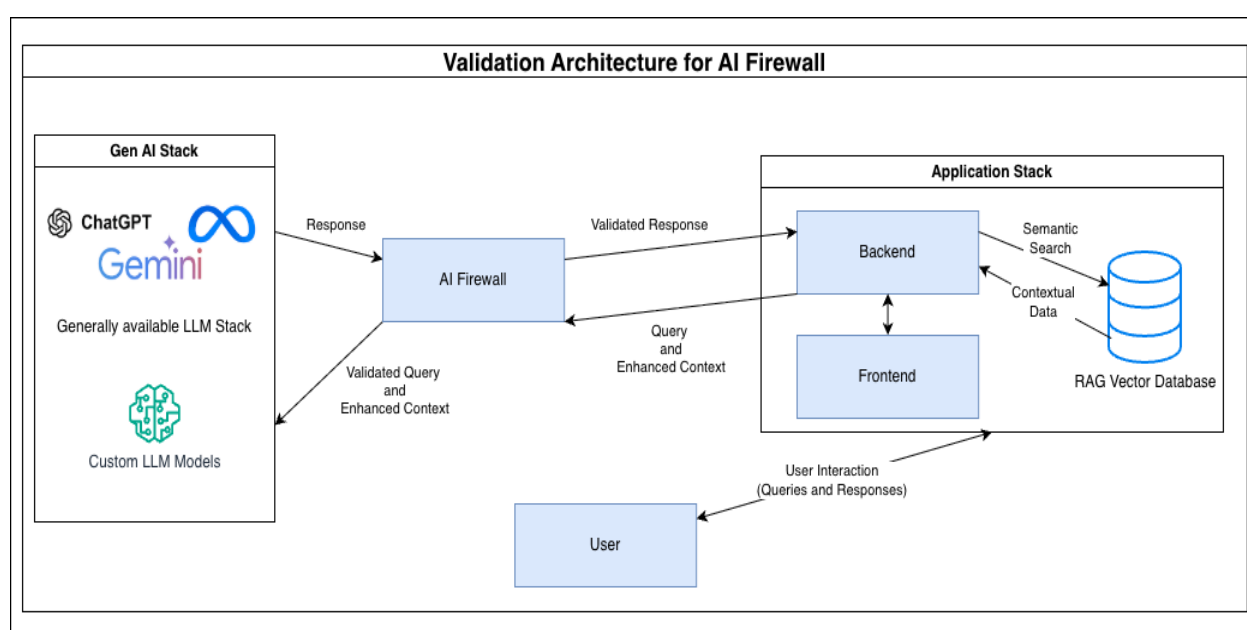


Figure 2: Validation Architecture for AI Firewall

The Validation Architecture includes validation steps at the AI Firewall layer:

- **Step 1: Outbound Validation (Prompt & Context Filtering)**
 - From Backend to Firewall: Instead of sending data directly to the LLM, the Backend sends the "Tokens and Query + Enhanced Context" to the AI Security.
 - The Inspection: The firewall scans this payload. It typically checks for information such as:
 - PII/Data Leakage: Ensuring no sensitive user data (names, credit cards) or confidential context is accidentally leaked.
 - Prompt Injection Attacks: Detecting malicious inputs designed to manipulate the LLM.
 - From Firewall to LLM Stack: Only if the data passes these checks does the firewall forward the

"Validated Tokens and Query + Enhanced Context" to the LLM Stack.

- **Step 2: Inbound Validation (Response Filtering)**
 - From LLM Stack to Firewall: The Response generated by the model is sent back to the AI Security solution first, rather than the backend.
 - The Inspection: The firewall scans the model's output. It typically checks information such as:
 - PII and Sensitive information Leakage
 - Hallucinations/Toxicity: Ensuring the response is safe, on-topic, and non-toxic.
 - Malicious Code: Preventing the LLM from generating harmful executable code.
 - From Firewall to Backend: The firewall forwards the "Filtered Response" to the Backend, which then safely displays it to the user via the Frontend.

6 SCORING MECHANISM FOR CYBERRISK VALIDATION

The scoring mechanism provides a structured approach to assess the effectiveness, usability, and compliance readiness of AI Security. Scores are based on the dual criteria of prevention (ability to stop threats/actions) and detection (ability to log, monitor, and audit activities).

Each efficacy dimension Security and Operational has tailored scoring criteria to ensure a holistic evaluation.

6.1 SECURITY EFFICACY SCORING CRITERIA

The Security Efficacy dimension measures how effectively the AI Security Control Layer detects, prevents, and logs cyber threats that directly impact the GenAI Application. Security efficacy is evaluated using a Prevention + Detection framework, where both active enforcement and audit visibility are considered essential security outcomes.

The effectiveness of security efficacy controls is graded as follows:

Outcome	Description	Score
Prevent and Detect	The threat or unauthorized action is actively blocked, and the event is logged with sufficient metadata to support auditability and forensic analysis.	100%
Prevent, No Detect	The threat is blocked, but no actionable logging or telemetry is generated, limiting audit and compliance visibility.	75%
Detect, No Prevent	The action is allowed to proceed, but the event is logged, providing post-event visibility without active mitigation.	25%
No Detect, No Prevent	The threat is neither blocked nor logged, representing a complete security control failure.	0%

6.2 OPERATIONAL EFFICIENCY SCORING CRITERIA

Deployment & Onboarding

Outcome	Description	Score
Seamless & Scalable	Deployment supports multiple models (proxy, SDK, gateway, sidecar), works across cloud, hybrid, and on-prem environments, requires minimal code changes, provides secure-by-default policies, and reaches enforcement rapidly with clear documentation and SSO/MFA support.	100%
Functional but Manual	Deployment supports common enterprise environments and LLM architectures but requires moderate manual configuration or code changes; onboarding and policy activation are achievable but not fully automated.	75%
Limited & Complex	Deployment works only for limited architectures or environments; onboarding requires significant manual effort, custom scripting, or extensive tuning before enforcement.	25%
Unsupported & Fragmented	Deployment is ad-hoc, lacks enterprise compatibility, fails to enforce policies reliably, or is unsuitable for production use.	0%

Policy Management & Administration

Outcome	Description	Score
Granular & Explainable	Supports rule-based and/or ML-assisted policies with fine-grained controls (user, app, data, intent, retrieval), clear explainability, versioning, rollback, audit history, and separation of duties.	100%
Configurable but Limited	Policies are configurable and enforceable but lack full granularity, clear explainability, or mature change management features.	75%
Coarse & Opaque	Policies are static or coarse-grained, explanations are unclear, and versioning or rollback is limited or manual.	25%
Hard-Coded or Unmanageable	Policies are fixed, undocumented, or cannot be safely changed without disrupting production.	0%

Integration with Enterprise Ecosystem

Outcome	Description	Score
Deeply Integrated	Native integrations with SIEM/SOAR/XDR, IAM/SSO, and DLP; identity-aware enforcement; robust APIs and real-time telemetry export.	100%

Partially Integrated	Integrates with common enterprise tools but requires custom connectors, limited telemetry, or delayed exports.	75%
Minimal Integration	Provides basic logs or APIs but lacks real-time integration or identity context.	25%
Isolated	No meaningful enterprise integration; operates as a standalone tool with limited visibility.	0%

Incident Response & Visibility

Outcome	Description	Score
SOC-Ready	Real-time alerts, full incident context (prompt, retrieval, output, action), strong correlation with enterprise telemetry, and forensically sound evidence retention.	100%
Detective but Reactive	Alerts and logs exist but lack full context or require manual correlation for investigations.	75%
Limited Visibility	Events are logged but poorly contextualized; alerts are delayed or noisy.	25%
Operationally Blind	Minimal alerting or logging; incidents cannot be reconstructed reliably.	0%

Insight for Threat Hunting & Forensics

Outcome	Description	Score
Proactive & Analytical	Supports prompt/response traceability, conversation timelines, attack pattern analytics, and advanced search/filtering across users, models, and data types.	100%
Reactive Visibility	Enables basic traceability and manual investigations but lacks advanced analytics or pattern detection.	75%
Event-Level Only	Provides isolated logs without session context or historical analysis capabilities.	25%
No Forensic Value	Insufficient data for threat hunting or post-incident analysis.	0%

Security Administration

Outcome	Description	Score
Enterprise-Grade Governance	Strong RBAC, multi-tenant isolation, immutable audit logs, compliance-ready reporting (OWASP, NIST), and operational resilience controls.	100%
Adequate Governance	RBAC and audit logs exist but multi-tenant or compliance reporting capabilities are limited.	75%

Basic Administration	Limited roles, partial logging, and manual compliance reporting.	25%
Operationally Risky	Weak or absent administrative controls; no auditability or governance support.	0%

7 GENERAL EVALUATION APPROACH

7.1 AI SECURITY VENDOR PARTICIPATION SELECTION CRITERIA

SecureIQLab will identify and invite vendors to participate in this AI Security CyberRisk Validation based on objective, transparent, and market-relevant criteria. Vendor selection is intended to ensure that the evaluation reflects a representative cross-section of commercially relevant AI Security solutions.

Market Leaders

Vendors that demonstrate significant market presence, as evidenced by:

- Material revenue attributable to AI Security or GenAI security offerings
- Broad enterprise deployment across multiple industries or geographies
- Established channel, cloud, or platform partnerships

Analyst-Recognized and Enterprise Challengers

Vendors identified through:

- Coverage or reference in independent industry analyst research, market reports, or buyer guides
- Input from enterprise security professionals, including CISOs, SOC teams, and platform owners
- Feedback obtained through surveys, direct inquiries, and interactions with enterprises, MSPs, and MSSPs

Innovative New Entrants

Vendors that:

- Offer differentiated or novel approaches to AI security
- Address risks identified in the OWASP Top 10 for LLM Applications
- Express interest in participating in independent testing to demonstrate technical capability

Conflict of Interest Statement

SecureIQLab affirms that it has no known conflicts of interest related to the selection of vendors for this validation. Vendor inclusion or exclusion is determined solely by the criteria outlined above. Participation does not constitute endorsement.

7.2 LIST OF CONSIDERED AI SECURITY VENDORS

The following AI Security solution have been considered for inclusion in this phase of validation:

Vendor Name	Product Name
SentinelOne	Prompt Security
Lasso Security	Lasso for Application

Imperva	Imperva AI Application Security
Veeam	Securiti AI
Radware	Radware LLM Firewall
Fortinet	FortiAI Gate
F5	AI runtime security
Check Point	Lakera
Netskope One	Netskope One
DeepKeep	DeepKeep LLM
Palo Alto Networks	Prisma AIRS
Akamai	Firewall for AI
Cloudflare	Firewall for AI
TrojAI	TrojAI Defend
Cato Networks	Cato AI Firewall
Hidden Layer	AISEC Platform
Google	Model Armor
Pillar Security	Pillar Platform

Note: Final participation is subject to availability, product readiness, and scope alignment.

7.3 VALIDATION TIMELINE

Schedule Timeline for Cyber Risk Validation

Index	Test Activity	Date Range	Dependencies
1	Test Commencement	04/01/2026	Vendor voluntary participation or procurement of vendor software
2	Confirm Vendor Configuration Feedback	05/01/2026	All vendors installed, smoke tested, and configurations validated
3	Testing	05/18/2026	Comprehensive testing
4	Feedback and Dispute Resolution / Retests	06/22/2026	Vendor feedback and dispute submissions

5	Publish Results	07/20/2026	Responsible disclosure requirements
---	-----------------	------------	-------------------------------------

Dates may be adjusted to accommodate vendor availability, remediation activities, or disclosure considerations

7.4 RISK AND RISK MANAGEMENT

At the time of publication, SecureIQLab is not aware of any material risks that would prevent execution of this test in accordance with this methodology. SecureIQLab will monitor and manage risks throughout the testing lifecycle.

7.5 GEO LIMITATION

SecureIQLab will make reasonable efforts to ensure that test traffic, adversarial simulations, and attack scenarios are not unduly biased toward a specific geographic location, except where necessary to validate relevant controls. Where applicable, traffic may originate from multiple regions to reflect global enterprise usage.

7.6 DISTRIBUTION OF TEST DATA

Upon completion of the validation phases:

- Individual vendor reports and one comparative report will be produced
- Results will be available to vendors and may be licensed for marketing purposes after the purchase of marketing rights.
- Public versions of reports will be available at <https://secureiqlab.com/publications/>

SecureIQLab retains editorial control over public disclosures.

7.7 FUNDING AGREEMENT

This is a non-commissioned, independently funded test conducted by SecureIQLab. Vendor participation does not influence scoring, methodology, or publication outcomes.

7.8 DISPUTE PROCESS

SecureIQLab will make reasonable efforts to resolve disputes related to test results or scoring

- Vendors receive detailed individual scorecards covering:
 - Security efficacy
 - Operational efficiency
- Each participating vendor will receive its individual results and scoring breakdown
- Vendors may submit disputes within two (2) weeks of receiving results
- Any validated scoring adjustments will be applied consistently across all vendors
- Vulnerabilities identified during testing will be handled under responsible disclosure practices
- SecureIQLab will not modify results after publication of final reports

This evaluation is designed to be fair, transparent, repeatable, and vendor-neutral, and is intended to comply with the principles of the AMTSO Testing Standards, including disclosure, comparability, and responsible reporting.

8 OPT-OUT POLICY

Opt-Out Eligibility

SecureIQLab will consider opt-out requests only under the following circumstances:

A. Out-of-Scope Determination

The product, solution, or technology is determined by SecureIQLab to be outside the defined scope of this AI Security Cyber Risk Validation methodology.

B. Lack of General Availability

The product, solution, or technology is not generally available, is in beta, preview, or otherwise not ready for enterprise deployment at the time of testing.

C. Public Interest Consideration

SecureIQLab determines that publication of test results would not serve the public interest.

Opt-out requests submitted for reasons outside those listed above will not be considered.

Opt-Out Submission Requirements

All opt-out requests must be submitted in writing.

Email submissions:

info@secureiqlab.com

Mailed submissions:

SecureIQLab LLC.
9600 Great Hills Trail Suite 150W
Austin, TX 78759

Required Information

To be considered valid, an opt-out request must include:

- Vendor name
- Product name
- Name and title of the authorized representative
- Email address and phone number
- The specific opt-out reason (from the eligible criteria above)
- Supporting details justifying the request

Incomplete requests will be considered invalid.

Opt-Out Timing

- The opt-out period begins at Test Commencement
- The opt-out period ends at the conclusion of the Dispute Phase

SecureIQLab will acknowledge receipt of an opt-out request and contact the vendor within three (3) business days to discuss feasibility and status.

Opt-Out Outcomes

- Opt-out before completion of the Configuration Phase:
The vendor will be listed in the results as
“Participant, not tested.”
- Opt-out after testing has been performed:
The vendor will be listed in the results as
“Tested, not published.”

In all cases, SecureIQLab retains discretion over classification wording to ensure consistency and transparency.

9 ATTESTATIONS

I understand and agree that I am submitting this Test Plan, and the following Attestations, on behalf of the entity listed below, and I represent and warrant that I have authority to bind such entities to these Attestations. All references to “I” or “me” or similar language refer to such an entity. I represent and warrant that the following Attestations are true, to the best of my knowledge and belief, and each of the following commitments will be upheld to the best of my ability.

I will provide public notification on the AMTSO website covering my obligation for notification of a Public Test, regardless of whether a potential Participant is in actual receipt of such notification prior to the Commencement Date of a Test.

All products included in this Test will be analyzed fairly and equally.

I will disclose any anticipated or known imbalance or inequity in the Test design to all Participants in the Test.

Although I may charge for participation in a Test, I will not charge any additional fees for a vendor to be a test subject under the Standards.

I will disclose any material conflicts of interest or other information that could materially impact the reliability of the Test.

I will disclose how the Test was funded.

I hereby affirm, to the best of my knowledge and belief that this Test Plan complies with the AMTSO Testing Standards, as of the date hereof.

Signature: /s/ David Ellis

Name: David Ellis

Test Lab: SecureIQLab

AMTSO Test ID: AMTSO-LS1-TP158

10 DOCUMENT AND REVISION

Version	Section	Revision overview
1.0	All	This is the first version of the methodology and the release candidate.

11 COPYRIGHT AND DISCLAIMER

Copyright © 2026 SecureIQLab, LLC. All rights reserved. The content of this report is protected by United States and international copyright laws and treaties. You may only use this report for your personal, non-commercial, informational purposes. Without SecureIQLab's prior written consent, you may not: (i) reproduce, modify, adapt, create derivative works from, publicly perform, publicly display, or distribute this report; or (ii) use this report, the SecureIQLab name, or any SecureIQLab trademark or logo as part of any marketing, promotion or sales activities. THIS REPORT IS PROVIDED "AS IS," "AS AVAILABLE" AND "WITH ALL FAULTS." TO THE MAXIMUM EXTENT PERMITTED BY LAW, SECUREIQLAB EXPRESSLY DISCLAIMS ALL WARRANTIES AND REPRESENTATIONS, EXPRESS OR IMPLIED, INCLUDING: (a) THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE; AND (b) ANY WARRANTY WITH RESPECT TO THE QUALITY, ACCURACY, CURRENCY OR COMPLETENESS OF THE REPORT, OR THAT USE OF THE REPORT WILL BE ERROR-FREE, UNINTERRUPTED, FREE FROM OTHER FAILURES OR WILL MEET YOUR REQUIREMENTS. WITHOUT LIMITING THE GENERALITY OF THE FOREGOING SENTENCE, YOU ACKNOWLEDGE AND AGREE THAT THE QUALITY, ACCURACY, CURRENCY AND COMPLETENESS OF THE REPORT DEPEND UPON VARIOUS FACTORS, INCLUDING FACTORS OUTSIDE OF SECUREIQLAB'S CONTROL, SUCH AS: (1) THE QUALITY, ACCURACY, CURRENCY OR COMPLETENESS OF INFORMATION AND MATERIALS PROVIDED BY OTHER PARTIES THAT ARE RELIED UPON BY SECUREIQLAB IN PERFORMING AND PREPARING THE REPORT; AND (2) THE UNDERLYING ASSUMPTIONS MADE BY SECUREIQLAB IN PREPARING THE REPORT REMAINING TRUE AND ACCURATE. YOU ARE SOLELY RESPONSIBLE FOR INDEPENDENTLY ASSESSING THE QUALITY, ACCURACY, CURRENCY AND COMPLETENESS OF THE REPORT BEFORE TAKING OR OMITTING ANY ACTION BASED UPON THE REPORT. IN NO EVENT WILL SECUREIQLAB BE LIABLE FOR ANY LOST PROFITS OR COST OF COVER, OR DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, PUNITIVE OR CONSEQUENTIAL DAMAGES, INCLUDING DAMAGES ARISING FROM OR RELATING TO ANY TYPE OR MANNER OF COMMERCIAL, BUSINESS OR FINANCIAL LOSS, EVEN IF SECUREIQLAB HAD ACTUAL OR CONSTRUCTIVE KNOWLEDGE OF THE POSSIBILITY OF SUCH DAMAGES AND REGARDLESS OF WHETHER SUCH DAMAGES WERE FORESEEABLE.

For more information about SecureIQLab and the testing methodologies, please visit our website. www.secureiqlab.com.

SecureIQLab (March 2026)

12 APPENDICES

12.1 VALIDATION FRAMEWORK MAPPING TABLE

This table provides a comprehensive mapping of each test case in the framework to its corresponding risk categories in the OWASP Top 10 for LLM Applications and the relevant adversary tactics and techniques from the MITRE ATLAS framework. It serves as a "Rosetta Stone" for GenAI security, translating specific technical validation activities into the widely understood languages of risk management and threat intelligence.

Part I – Input Security (TC1–TC8)

TC #	Test Case Name	OWASP LLM Top 10 Mapping	MITRE ATLAS Tactic(s)	MITRE ATLAS Technique(s)
TC1	Direct Prompt Injection	LLM01:2025 Prompt Injection	Initial Access, Defense Evasion	LLM Prompt Injection
TC2	Indirect Prompt Injection	LLM01:2025 Prompt Injection	Initial Access, Defense Evasion	LLM Prompt Injection
TC3	Multilingual / Obfuscated Injection	LLM01:2025 Prompt Injection	Defense Evasion	Obfuscated Prompt Injection
TC4	Vector Database Toxic Content Injection	LLM03: Training Data Poisoning	Resource Development, Impact	Poison Training Data
TC5	PII Disclosure via Direct Input	LLM02: Sensitive Information Disclosure	Collection	LLM Data Leakage
TC6	Payment Card Data Ingestion (PCI)	LLM02: Sensitive Information Disclosure	Collection	LLM Data Leakage
TC7	Oversized Prompt Blocking	LLM10:2025 Unbounded Consumption	Impact	Denial of ML Service
TC8	High-Volume Request Flood	LLM04: Model Denial of Service	Impact	Denial of ML Service, Cost Harvesting