

# AI, LLM, and Agent Security Methodology and Testing Checklist

---

<b>Document Version</b>	1.0
<b>Prepared By</b>	Lorikeet Security
<b>Contact</b>	lorikeetsecurity.com



# 1. Executive Overview

*A specialized assessment framework for AI systems including large language models, autonomous agents, RAG pipelines, chat interfaces, and AI-powered applications.*

Lorikeet Security's AI and LLM penetration testing practice addresses the unique attack surface of modern AI systems. This includes prompt injection and jailbreaking, insecure output handling, agentic tool use exploitation, model supply chain risks, training data leakage, and safety control bypass. Our methodology is aligned to the OWASP LLM Top 10 (2025), NIST AI RMF, and MITRE ATLAS adversarial threat framework.

AI security assessments require specialized expertise distinct from traditional application testing. Our practitioners have hands-on experience with LLM architectures, agentic frameworks (LangChain, AutoGPT, CrewAI, custom tool-calling pipelines), RAG systems, and both open-source and commercial model deployments. We assess AI systems as both standalone targets and as components within broader application architectures.

## 2. Assessment Methodology

*Three-phase model: System Mapping, Active Red-Teaming, and Risk Reporting.*

### Phase 1: System Mapping and Threat Modeling

Prior to active testing, Lorikeet Security maps the complete AI system architecture including model provider, deployment configuration, tool/function integrations, data pipelines, RAG knowledge bases, and user-facing interfaces. A threat model is constructed identifying trust boundaries, data flows, and highest-risk attack surfaces based on the system's capabilities and access to sensitive resources.

### Phase 2: Active Red-Teaming

Active testing combines manual adversarial prompting by security researchers with systematic coverage across all nine assessment domains. Testing is calibrated to the model's deployment context: a customer-facing chatbot with limited tool access receives different testing depth than an autonomous agent with access to email, databases, and code execution.

**Black-box:** Black-box: no system prompt, model card, or architecture access

- Simulates an external attacker or malicious end user

**Gray-box:** Gray-box: system prompt, tool list, and deployment context provided

- Simulates an insider threat or disclosed-architecture adversary

**White-box:** White-box: full architecture, system prompts, fine-tuning details, and pipeline code

- Maximum coverage; recommended for high-risk agentic deployments

## Phase 3: Risk Reporting and Remediation Guidance

Findings are documented with reproduction steps, attack narrative, business impact analysis, and specific remediation guidance tailored to the AI system's architecture. Recommendations are classified as architectural changes, prompt engineering controls, infrastructure hardening, or monitoring improvements.

### 3. Scope of Assessments

*Standard scope covers the AI system's user interface, inference API, tool integrations, data pipelines, and surrounding application infrastructure.*

Component	Scope Details
<b>LLM / Foundation Model</b>	Prompt injection resistance, jailbreaking, output safety, and behavioral boundary testing
<b>Chat Interface</b>	Web/mobile UI, conversation history, file handling, sharing features, and client-side security
<b>Agentic Tool Integrations</b>	All tools, function calls, MCP servers, and external APIs accessible to the agent
<b>RAG Pipeline</b>	Knowledge base ingestion, retrieval mechanisms, embedding storage, and corpus poisoning vectors
<b>Inference API</b>	Authentication, authorization, rate limiting, model serving infrastructure security
<b>Data Pipelines</b>	Training data ingestion, fine-tuning datasets, and corpus update mechanisms
<b>Memory Systems</b>	Short-term context, long-term memory stores, and cross-session persistence mechanisms
<b>Multi-Agent Orchestration</b>	Trust boundaries between agents, inter-agent communication security, and escalation paths

### 4. Tools and Techniques

*AI security testing combines adversarial ML research techniques with traditional security tooling adapted for AI system attack surfaces.*

Domain	Tools	Techniques
<b>Prompt Injection</b>	Garak, custom prompt libraries, LLM Red Teaming frameworks	Direct/indirect injection, multi-turn escalation, encoding bypass
<b>Jailbreaking</b>	Custom adversarial prompts, many-shot testing, GCG-style suffix generation	Persona manipulation, roleplay bypass, cross-lingual evasion
<b>Agent Testing</b>	Custom tool harnesses, MCP client tooling, LangChain test environments	Tool abuse, privilege escalation, trust boundary violation

Domain	Tools	Techniques
Data Extraction	Membership inference prompts, extraction probes, embedding analysis	Training data leakage, PII recovery, system prompt extraction
Output Analysis	Burp Suite, custom middleware, LLM output parsers	XSS via output, SSRF generation, downstream injection testing
Infrastructure	Nmap, nuclei, API testing tools, model serving CVE scanners	Inference API hardening, network isolation, credential exposure
RAG Testing	Custom corpus injection scripts, vector DB query analysis	Poisoning, retrieval manipulation, cross-user data leakage

## 5. Testing Checklist by Domain

*Comprehensive test cases across nine assessment domains covering the full OWASP LLM Top 10 and agentic attack surface.*

Risk Level Key: Critical = Immediate remediation required | High = Remediate within 30 days | Medium = Remediate within 90 days | Low = Informational

ID	Test Case	Standard	Risk Level
<b>1. Prompt Injection and Jailbreaking</b>			
PI-01	Test for direct prompt injection: inject adversarial instructions via user-controlled input fields to override system prompt behavior	OWASP LLM01	Critical
PI-02	Test for indirect prompt injection via external data sources: documents, emails, web pages, RAG corpus, and tool outputs consumed by the agent	OWASP LLM01	Critical
PI-03	Assess jailbreaking resistance: attempt role-play personas, hypothetical framing, token smuggling, and encoding tricks to bypass safety guardrails	OWASP LLM01	High
PI-04	Test for system prompt extraction: craft inputs designed to leak the system prompt, pre-prompt instructions, or operational context to the user	OWASP LLM07	High
PI-05	Evaluate multi-turn injection persistence: inject instructions in early turns and test whether they persist or escalate across a conversation session	OWASP LLM01	Critical
PI-06	Test for cross-context injection in multi-modal inputs: images with embedded text instructions, audio transcription injection, and PDF/document injection	OWASP LLM01	High
PI-07	Assess injection resistance in agentic pipelines where LLM output is passed as input to a downstream LLM or tool without sanitization	OWASP LLM08	Critical
<b>2. Insecure Output Handling</b>			
OUT-01	Test for LLM output passed directly to downstream systems: shell execution, SQL queries, code interpreters, or API calls without sanitization	OWASP LLM02	Critical
OUT-02	Assess for XSS via LLM output rendered in web interfaces: craft inputs that cause the model to output HTML/JS that executes in the browser	OWASP LLM02	High

ID	Test Case	Standard	Risk Level
OUT-03	Test for SSRF via LLM-generated URLs: induce the model to produce URLs pointing to internal infrastructure consumed by backend fetch operations	OWASP LLM02	Critical
OUT-04	Evaluate whether LLM-generated code is executed in a sandboxed environment; test sandbox escape vectors where code execution is permitted	OWASP LLM02	Critical
OUT-05	Assess markdown/LaTeX injection in rendering contexts: test for formula injection and hyperlink injection in platforms that render LLM output directly	OWASP LLM02	Medium
<b>3. Training Data and Model Supply Chain</b>			
TRAIN-01	Assess data poisoning risk: evaluate controls on training data ingestion pipelines, fine-tuning datasets, and RAG corpus update mechanisms	OWASP LLM03	Critical
TRAIN-02	Test for training data extraction: craft membership inference and extraction prompts to probe for memorized PII, credentials, or sensitive organizational data	OWASP LLM06	High
TRAIN-03	Review model supply chain integrity: verify checksums, provenance attestations, and access controls on base model artifacts and fine-tuned weights	OWASP LLM03	High
TRAIN-04	Assess RAG retrieval poisoning: inject adversarial documents into the knowledge base and test whether they alter model behavior or exfiltrate retrieved content	OWASP LLM01/03	Critical
TRAIN-05	Evaluate embedding model security: test for nearest-neighbor extraction attacks and embedding inversion to recover sensitive source documents	OWASP LLM03	Medium
<b>4. Excessive Agency and Tool Use</b>			
AGEN T-01	Assess tool/function call authorization: verify the agent cannot be induced to call tools outside its intended scope via prompt injection or role confusion	OWASP LLM08	Critical
AGEN T-02	Test for privilege escalation via tool chaining: induce the agent to chain tool calls in sequences that collectively exceed authorized permissions	OWASP LLM08	Critical
AGEN T-03	Evaluate least-privilege enforcement on agent tool permissions: verify each tool is scoped to the minimum required access (read-only vs. write, specific APIs)	OWASP LLM08 / NIST AI RMF	High
AGEN T-04	Test human-in-the-loop controls: verify that high-risk actions (sending emails, executing code, modifying data, making purchases) require explicit user confirmation	OWASP LLM08	Critical
AGEN T-05	Assess multi-agent trust boundaries: test whether a compromised or adversarially prompted sub-agent can issue instructions to a higher-privilege orchestrator agent	OWASP LLM08	Critical
AGEN T-06	Evaluate agent memory systems for injection persistence: test whether adversarial content stored in long-term agent memory influences future sessions	OWASP LLM08	High
AGEN T-07	Test computer use and browser agents for navigation hijacking: inject instructions via web content to redirect agent actions to attacker-controlled resources	OWASP LLM01/08	Critical
<b>5. Sensitive Information Disclosure</b>			
INFO-01	Test for PII leakage via context window: craft prompts that cause the model to repeat or summarize prior user inputs containing sensitive information	OWASP LLM06	High

ID	Test Case	Standard	Risk Level
INFO-02	Assess cross-user data leakage in multi-tenant deployments: test whether one user's context bleeds into another user's session via shared model state or caching	OWASP LLM06	Critical
INFO-03	Evaluate credential and API key exposure via tool responses: verify that secrets returned by tool calls are not echoed back in plaintext model responses	OWASP LLM06	Critical
INFO-04	Test for model fingerprinting and intellectual property disclosure: probe for proprietary fine-tuning data, internal documents, and confidential business logic encoded in the model	OWASP LLM06	Medium
INFO-05	Assess logging and observability pipeline for sensitive data retention: verify that prompt/response logs are redacted, access-controlled, and have defined retention limits	OWASP LLM06 / GDPR	High
<b>6. Denial of Service and Resource Exhaustion</b>			
DOS-01	Test for prompt-based resource exhaustion: submit inputs designed to maximize token generation, chain-of-thought depth, or recursive tool calls	OWASP LLM04	High
DOS-02	Assess rate limiting and cost controls on inference endpoints: verify per-user, per-session, and global token budgets are enforced	OWASP LLM04	High
DOS-03	Test for ReDoS (catastrophic backtracking) in pre/post-processing regex patterns applied to LLM inputs and outputs	OWASP LLM04	Medium
DOS-04	Evaluate agent loop detection: verify the system detects and terminates infinite tool-call loops or self-referential reasoning chains	OWASP LLM08	High
DOS-05	Assess context window exhaustion attacks: inject padding or verbosity-inducing prompts to push legitimate content out of the context window	OWASP LLM04	Medium
<b>7. Model and Inference API Security</b>			
API-01	Assess inference API authentication and authorization: verify API keys are scoped, rotated, and not exposed in client-side code or public repositories	OWASP LLM09	Critical
API-02	Test model serving infrastructure for known CVEs in frameworks (Triton, TorchServe, vLLM, Ollama) and verify network isolation of inference endpoints	OWASP LLM09	High
API-03	Evaluate model stealing risk via repeated inference calls: assess whether API responses can be used to distill a functional replica of the model	OWASP LLM10	Medium
API-04	Test system prompt confidentiality at the API level: verify that system prompts cannot be retrieved via API metadata endpoints or timing side-channels	OWASP LLM07	High
API-05	Assess streaming response handling for injection in partial completions: verify security controls apply to streamed tokens, not only complete responses	OWASP LLM02	Medium
<b>8. Guardrail and Safety Control Bypass</b>			
GUAR D-01	Test input classifier bypass: craft adversarial inputs using synonyms, leetspeak, Unicode homoglyphs, Base64 encoding, and multilingual obfuscation to evade content filters	OWASP LLM01	High
GUAR D-02	Assess output classifier bypass: produce outputs in formats that evade post-processing filters (structured data, code comments, steganographic output)	OWASP LLM02	High

ID	Test Case	Standard	Risk Level
<b>GUAR D-03</b>	Test for many-shot and few-shot jailbreaking: provide in-context examples that gradually normalize prohibited behavior and test compliance boundary erosion	OWASP LLM01	<b>High</b>
<b>GUAR D-04</b>	Evaluate persona and roleplay bypass resistance: test whether assigning the model a fictional role or character causes it to abandon safety behaviors	OWASP LLM01	<b>High</b>
<b>GUAR D-05</b>	Assess cross-lingual safety: test guardrail coverage for low-resource languages and mixed-language inputs where safety training may be underrepresented	OWASP LLM01	<b>Medium</b>
<b>GUAR D-06</b>	Test for token-level adversarial suffixes (GCG-style): evaluate whether adversarially optimized token sequences can reliably elicit policy violations	OWASP LLM01	<b>High</b>
<b>9. Chat Interface and UX Security</b>			
<b>UX-01</b>	Test chat interface for stored XSS via LLM-generated content rendered in the conversation UI	OWASP LLM02 / WSTG	<b>High</b>
<b>UX-02</b>	Assess conversation history handling: verify history is isolated per user, encrypted at rest, and not accessible to other users or session contexts	OWASP LLM06	<b>High</b>
<b>UX-03</b>	Test for CSRF on chat API endpoints and conversation management actions (delete, share, export)	OWASP WSTG-SESS-05	<b>Medium</b>
<b>UX-04</b>	Evaluate conversation sharing functionality for authorization bypass: test whether shared conversation links expose system prompts or other users' data	OWASP LLM06	<b>High</b>
<b>UX-05</b>	Assess file and attachment handling in multimodal chat: test for malicious file processing, SSRF via image URLs, and path traversal in document uploads	OWASP WSTG-BUSL	<b>High</b>

## 6. Compliance Standards Reference

*Lorikeet Security aligns AI and LLM assessments to the following primary standards and frameworks.*

Standard	Relevance
<b>OWASP LLM Top 10 (2025)</b>	Primary reference framework for LLM and AI application security; all findings are mapped to applicable OWASP LLM categories.
<b>NIST AI Risk Management Framework (AI RMF)</b>	Governs risk identification, measurement, and management for AI systems; applied during agentic system and model supply chain assessment.
<b>NIST SP 800-53 (AI Controls)</b>	Supplemental security controls applied to AI/ML system components including access control, audit logging, and system integrity.
<b>MITRE ATLAS</b>	Adversarial Threat Landscape for AI Systems; used as a reference for attack technique coverage mapping across all assessment domains.
<b>EU AI Act (Risk Classification)</b>	Applied where the customer operates in EU-regulated contexts; informs risk classification of high-risk AI system components.
<b>ISO/IEC 42001</b>	AI management system standard; applied during compliance-oriented assessments requiring documented AI governance controls.

Standard	Relevance
GDPR / CCPA	Privacy regulation requirements applied during sensitive data disclosure and logging assessment components.
SOC 2 Type II	Applied for SaaS AI product assessments where customer data processing and confidentiality controls are in scope.

## 7. Deliverables and Engagement Model

*Standard deliverables for every AI and LLM security assessment engagement.*

- **Final Report:** Final report: executive summary, OWASP LLM-mapped findings with CVSS scoring, reproduction prompts, and remediation roadmap
- **Prompt Library:** Adversarial prompt library: documented set of successful and near-miss attack prompts for regression testing during development
- **Architecture Review:** Architecture risk assessment: evaluation of system design decisions contributing to identified vulnerabilities with recommended architectural alternatives
- **Evidence Package:** Raw evidence package: prompt/response logs, tool call traces, and supporting screenshots organized by finding
- **Debrief Session:** Debrief session with AI/ML engineering and security teams; red team walkthrough of most impactful attack chains

Typical engagement timelines: 3-5 business days for a single chatbot or API endpoint; 7-14 business days for complex agentic deployments with multiple tool integrations, RAG pipelines, and multi-agent orchestration.