

AI SECURITY ENGINEERING FIELD GUIDE · 2026

# Chapter 14 · Secure AI Architecture Design

Standalone learning module for LMS delivery and required reading.

FORMAT

**Standalone PDF**

USE

**LMS reading  
module**

SCOPE

**Single chapter**

AUDIENCE

**Learners**

---

CHAPTER 14

# Secure AI Architecture Design

# Design

SECURE AI ARCHITECTURE

Put enforcement where the model can be wrong and the system still stays safe.

Put enforcement where the model can be wrong and the system still stays safe.

FIELD GUIDE

#### FIELD GUIDE: ARCHITECTURE REVIEW

Inspect trust placement, data-plane controls, model-provider boundaries, retrieval paths, tool authority, fallback routes, telemetry, and release gates.

This domain covers secure architecture for LLM applications, RAG systems, agents, workflows, model chains, provider boundaries, tool integrations, telemetry, and governance evidence. It matters when one design decision can determine whether prompt injection, retrieval leakage, excessive agency, privacy exposure, or missing evidence becomes a systemic failure.

Architecture review is where local fixes become a system design. Prompts, filters, evals, retrieval controls, tool policy, provider boundaries, approvals, and telemetry should not all fail for the same reason or rely on the same model judgment.

#### FIELD USE

Draw the architecture with three overlays: data flow, authority flow, and evidence flow. A design is not review-ready until all three are visible.

# AI VENDOR PROCUREMENT MAP

Buy trustworthy AI. Manage risk. Drive value. Ensure compliance.

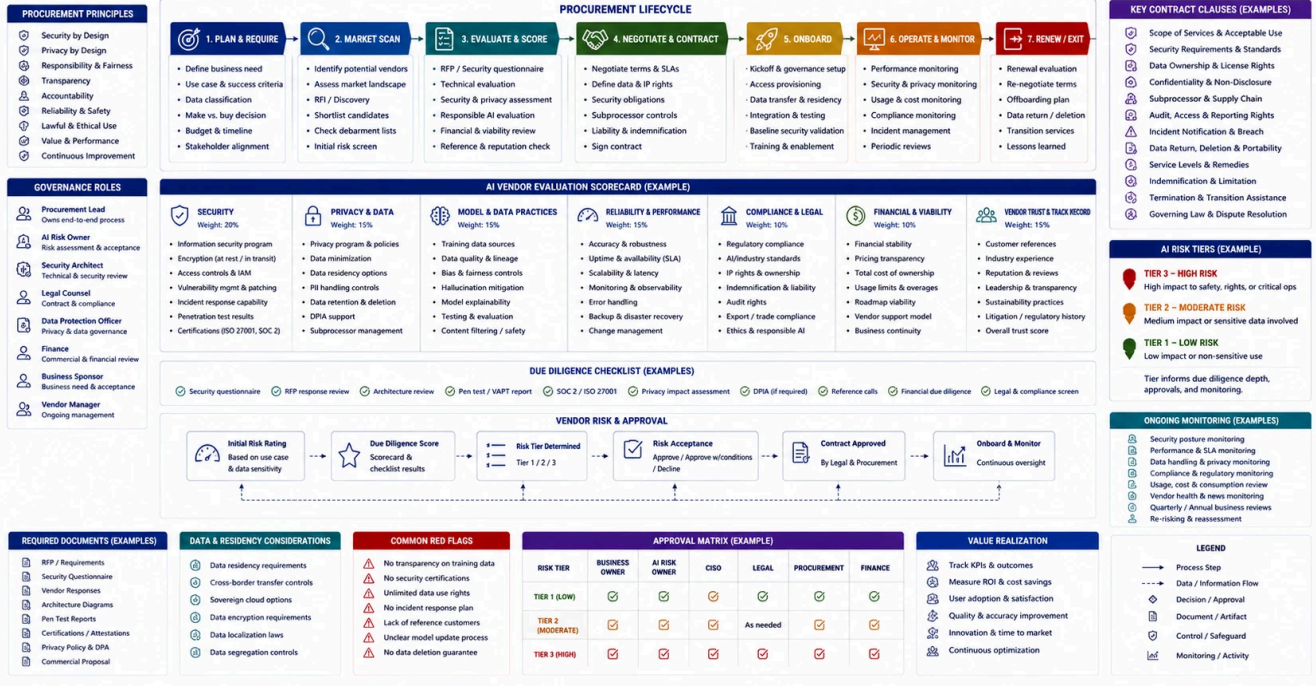


FIGURE 1: FIGURE 15: SECURE AI ARCHITECTURE REFERENCE SHOWING DATA FLOW, AUTHORITY FLOW, AND EVIDENCE FLOW OVERLAYS ACROSS PROMPT, RETRIEVAL, TOOL, PROVIDER, AND LOGGING LAYERS

## WHAT BREAKS HERE?

- The model becomes the sole enforcement point for authorization, policy, or safety.
- Retrieved content, tool output, and memory share one trust level.
- Fallback paths bypass security controls.
- Multi-model chains lose provenance and accountability.
- Architecture diagrams omit logs, approvals, exception paths, and control owners.
- Defense layers all depend on the same compromised context or model judgment.



## WHAT TO TEST

- › Model-enforced authorization failure.
- › Prompt injection across retrieval, tools, memory, and multi-model chains.
- › Fallback bypass of logging, authorization, rate limits, or approvals.
- › Agent blast-radius scenarios.
- › Cross-tenant, cross-role, and cross-provider data flows.
- › Evidence reconstruction from architecture trace points.
- › Control failure combinations across prompt, retrieval, tool, provider, and approval layers.

## CONTROLS AND GUARDRAILS

- › Explicit trust model for users, prompts, context, retrieval, tools, memory, providers, and logs.
- › Deterministic controls for authorization, data access, tool policy, approvals, and output handling.
- › Defense layers with independent failure modes.
- › Fallback security requirements and control-preserving degradation paths.
- › Architecture decision records for provider, retrieval, agent, data, and telemetry choices.
- › Release gates tied to architecture risks and evidence production.
- › Independent layers for data access, action authorization, output handling, telemetry, and approvals.



## RELATED SERVICES AND WORKBENCH TOOLS

TYPE	RELATED PATHS
Workbench	Threat Canvas, Authority Graph, Runtime Proxy, AI Control Crosswalk
Services	AI Product Security Assessment, Product Security Baseline, Agentic Workflow Hardening
Handbook	<a href="#">Architecture and Trust Boundaries</a> , <a href="#">Threat Modeling</a> , <a href="#">Governance Evidence and Customer Trust</a>

### DECISION - CONDITIONAL: ARCHITECTURE READINESS DECISION

Approve the design only when enforcement, evidence, fallback behavior, and ownership remain clear even if model behavior is unreliable.