

Chapter 03 · Prompt Injection and Context Security

Standalone learning module for LMS delivery and required reading.

FORMAT

Standalone PDF

USE

**LMS reading
module**

SCOPE

Single chapter

AUDIENCE

Learners

Prompt Injection and Context Security

INSPECT

System, developer, user, retrieved, tool, and memory context boundaries.

TEST

Direct injection, indirect injection, tool-output injection, memory persistence, and policy bypass.

Treat untrusted context as data. Never let it become authority.

FIELD GUIDE

FIELD GUIDE: CONTEXT REVIEW

Inspect every place instructions, data, retrieval results, tool output, memory, and conversation history meet inside the model context.

This domain covers direct prompt injection, indirect prompt injection, context confusion, instruction hierarchy, context provenance, memory contamination, tool-output injection, and prompt-driven policy bypass. It matters when an AI system reads untrusted content from users, web pages, documents, tickets, repositories, emails, retrieval systems, browser sessions, or tools.

The practical question is not whether the model can be tricked in isolation. The useful question is what a successful instruction conflict can reach. A low-impact summary error is different from a poisoned document that changes a tool call, exposes sensitive context, or persists in memory.

FIELD USE

Classify each context source by authority. Then test whether a lower-authority source can change a higher-authority decision, tool call, citation, memory entry, or data-handling rule.

CONTEXT AUTHORITY STACK

Separate data from authority. Lower-trust context cannot change higher-authority decisions.

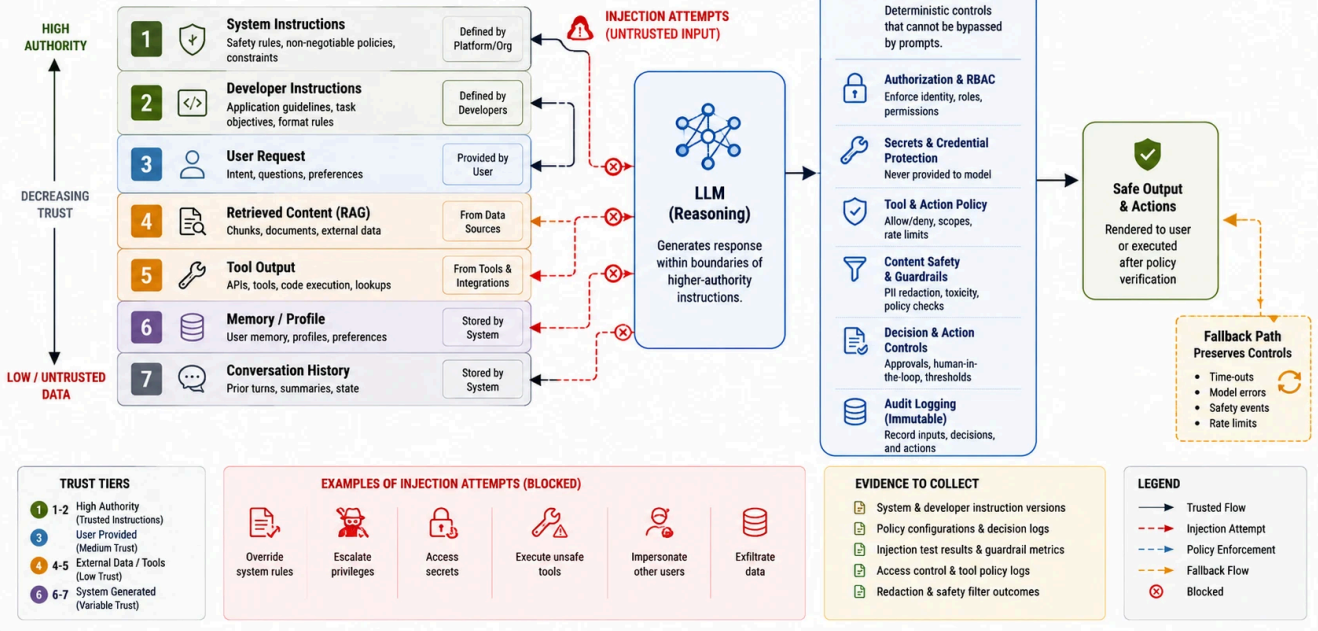


FIGURE 1: FIGURE 4: CONTEXT AUTHORITY STACK SHOWING TRUST LEVELS FROM SYSTEM INSTRUCTIONS THROUGH DEVELOPER, USER, RETRIEVED, TOOL, AND MEMORY SOURCES

WHAT BREAKS HERE?

- Untrusted content overrides system or developer instructions.
- Retrieved documents carry hidden instructions into context.
- Tool output is treated as policy, permission, or user intent.
- Memory stores poisoned content and replays it later.
- The application relies on the model to separate data from authority.
- A successful injection is dismissed because the final answer looked harmless.

EVIDENCE TO COLLECT

- › Context assembly diagram.
- › Prompt-injection test results.
- › Trust-tier table.
- › Memory policy and replay tests.
- › Tool-output handling notes.
- › Regression suite evidence.
- › Finding severity rationale.

OUTPUT ARTIFACTS

- › Context security review.
- › Prompt-injection findings register.
- › Context trust-tier matrix.
- › Regression test suite.
- › Remediation backlog.

RELATED SERVICES AND WORKBENCH TOOLS

TYPE	RELATED PATHS
Workbench	RAG Test Harness, Adversarial Range, Runtime Proxy
Services	AI Product Security Assessment, AI Red Team Validation, Product Security Baseline
Handbook	Prompt Injection , RAG Authorization , Agentic Permissions

DECISION - READY: CONTEXT AUTHORITY DECISION

If a source is not allowed to grant authority, the system must label it, constrain it, and prevent it from changing controls.

Prompt Injection and Context Security AISECURITY.LLC