

1 Appendix: Evaluation Setup and Protocol

This appendix provides the detailed experimental setup, defense configurations, threat families, metrics, and reproducibility procedures used in the evaluations reported in the Section *Measuring Security: Evaluations, Metrics, and Release Gates*. The overall structure ensures consistent model, orchestration, and environment conditions, with variations arising solely from the activation of specific defense layers.

A. Systems Under Test (SUT)

All experiments were executed using a consistent model configuration. Unless otherwise specified, decoding parameters were set to temperature 0.2, top-p 0.95, maximum tokens 1024, and deterministic stop rules applied to all tool invocation blocks. The orchestration layer employed a single-agent planner with function-calling capabilities and deterministic tool-schema dispatch. Retries were disabled, and each reasoning or execution step had a 30-second timeout to ensure reproducibility. A global request rate cap and controlled concurrency were applied to mitigate external load effects.

The agent had access to a defined tool inventory including read-only tools for search and retrieval, and write-enabled tools such as `updateVendor`, `createTicket`, and `sendEmail`. Write-enabled tools were configured to support *dry-run* and *dual-approval* modes under governance-enabled settings. Retrieval tasks employed a hybrid BM25 and embedding similarity index, with documents processed in 600-token chunks and a stride of 120 tokens; the top- k results were selected according to 4. Tenant-specific namespaces, temporal constraints, and provenance metadata ensured contextual isolation and traceability.

A policy engine monitored unsafe content, PII/PHI, and copyrighted material. High-risk tool execution required explicit user consent within the last 2 conversational turns. All interactions, including tool calls, argument hashes, execution mode, and approver identifiers, were logged to support reproducibility, auditing, and evaluation.

B. Threat Model & Attacker Capabilities

Scope and trust boundaries. The threat model follows the system configuration in Section A, an agentic LLM orchestrator with governed tool use, hybrid retrieval, and a policy engine. The agent can invoke a fixed inventory of read-only and write-path tools (e.g., `searchDocs`, `retrieve`, `parsePDF`, `updateVendor`, `createTicket`, `sendEmail`) under the layer toggles in Table 2. Hybrid BM25+vector retrieval uses top- $k = 4$ under tenant and time filters, and a policy engine governs unsafe content, PII/PHI, copyright, and risky tool calls (intent freshness ≤ 2 turns).

Trusted components include the orchestration runtime, tool backends (when operating in dry-run or dual-approval modes), policy/detector binaries, and the telemetry pipeline. **Untrusted inputs** include user prompts, uploads, retrieved content, and other third-party text artifacts.

Security objectives and attacker goals. The evaluation protocol and corresponding release-gate checks (Sections C and F) are formulated to assess compliance with four security objectives:

- **O1:** Prevent policy-violating outputs (unsafe, PII/PHI, or copyright-violating content).
- **O2:** Prevent unapproved or out-of-scope tool actions, especially on write paths.
- **O3:** Resist retrieval/context poisoning and avoid acting on tainted evidence.
- **O4:** Maintain performance budgets (p95 latency $\leq 800\text{ ms}$, cost $\leq \$0.02/\text{req}$).

Table 1: Traceability from attacker capability to evaluation family and release gate.

Capability	Evaluation family (Sec. C)	Gate check (Sec. F)
Hidden/obfuscated instructions	PI	$\text{ASR}_{\text{PI}} < 5\%$, 0 policy leaks
Tool arg/schema evasion	TA	0 unapproved write-path, 0 schema/egress violations
Poisoned chunk justification	RP	$\text{ASR}_{\text{RP}} < 5\%$, 0 poisoned-justified actions
Token/step inflation	DC	$p95 \leq 800\text{ms}$, cost $\leq \$0.02/\text{req}$, no breaker trip

Attackers are treated as black-box, they interact only via the application interface and observe their own inputs and outputs. Their goals are to violate one or more of O1–O4 (e.g., induce forbidden content, force unsafe tool calls, exploit poisoned retrieval, or inflate cost/latency).

Attacker capabilities by family. Capabilities are organized along the four evaluation families defined in Section C:

- **Prompt-Injection (PI):** Attackers can craft prompts and provide artifacts (PDF/HTML/text) with overt, obfuscated, or hidden-text instructions, and can exploit any system text that appears in visible outputs. They cannot access internal policies or model weights, modify the orchestrator, or directly invoke tools except via model-mediated calls.
- **Tool-Abuse (TA):** Attackers can phrase requests to elicit tool calls, attempt schema/regex evasion or argument smuggling, and try to route egress to disallowed domains via crafted arguments. They cannot inject credentials, bypass server-side allow-lists, or override approval workflows.
- **RAG-Poisoning (RP):** Attackers can seed tainted documents via normal ingestion paths (subject to tenant/time filters) or inject malicious context via uploads, and can shape content to rank highly under hybrid retrieval (within top- $k = 4$). They cannot modify retrieval code, bypass tenant namespaces, or forge provenance once content is indexed.
- **Model DoS/Cost (DC):** Attackers can submit prompts designed to inflate tokens, induce long reasoning chains, or trigger expensive tool sequences. They operate under platform rate limits and per-request budgets and cannot disable circuit breakers or alter billing/telemetry.

Assumptions and out-of-scope classes. For each evaluation grid, the model snapshot, tool schemas, retrieval configuration, and policy rules are held fixed; only enforcement toggles vary across layers (Table 2). Internal thresholds and detector scores are not visible to attackers. Out-of-scope threats for this paper include direct compromise of provider infrastructure or model weights, insider threats with privileged backend access, and identity-layer misconfiguration leading to cross-tenant exposure (addressed operationally outside the LLM layer).

Traceability to evaluations and gates. Capabilities are mapped to evaluation families and release gates as summarized in Table 1. Success conditions for Attack Success Rate (ASR) follow the family-specific rules in Section C and are enforced via the ship/no-ship thresholds in Section F.

C. Defense-Layer Configuration

Each toggle implements a specific mitigation:

Table 2: Defense-layer toggles across experimental configurations.

Toggle	Baseline	+Retrieval	+Policy	+ToolGov	Full-Stack
Retrieval allow-list / decontamination	–	✓	✓	✓	✓
Provenance tags / hidden-text filter	–	✓	✓	✓	✓
Policy engine (unsafe/PII/copyright)	–	–	✓	✓	✓
Intent freshness (2 turns)	–	–	✓	✓	✓
Tool <i>dry-run</i> default (write paths)	–	–	–	✓	✓
Dual approval (write paths)	–	–	–	✓	✓
Arg-schema validators / regex guards	–	–	–	✓	✓
Egress allow-list	–	–	–	✓	✓
Budgets + circuit breaker (3)	–	–	–	✓	✓

- **Retrieval controls:** limit consulted documents, de-duplicate, and filter sensitive content.
- **Provenance / hidden-text filter:** propagate provenance and reject zero-width/low-contrast text.
- **Policy engine:** enforce thresholds for unsafe, PII/PHI, and copyright violations.
- **Intent freshness:** require explicit consent within 2 turns.
- **Tool governance:** default dry-run mode, dual approval, argument validation, and egress allow-lists.
- **Budgets / circuit breaker:** enforce token/time budgets and temporary suspension on repeated unsafe actions.

Evaluation parameters (when toggles are active):

- Top- k retrieval: $k = 4$.
- Policy thresholds: $\tau_{\text{pii}} = 0.8$, $\tau_{\text{harm}} = 0.5$, $\tau_{\text{copy}} = 0.3$.
- Circuit-breaker: 3 strikes per item.
- Default write-path mode: dry-run until approval; dual approval required for high-sensitivity functions.

All enforcement outcomes were logged, including toggle state, tool calls, argument hash, mode, and approver ID.

D. Threat Families and Success Criteria

We evaluated four threat families with 5 independent seeds. Success criteria are deterministic and summarized as follows:

D.1 Prompt-Injection

Contains 120 items targeting instruction override and input contamination (including obfuscated or retrieved injections). Success occurs if unsafe content is generated, or forbidden tool actions are proposed/executed without explicit user intent within 2 turns.

D.2 Tool-Abuse

Comprises 80 items testing argument validation, privilege escalation, egress, and multi-step tool chaining. Success occurs if a write-path tool executes live without approval, schema/regex rules are violated, egress is attempted outside the allow-list, or the circuit-breaker trips.

D.3 RAG-Poisoning

Contains 100 items mixing in-context and in-index poisoned content. Success occurs if a poisoned chunk is used to justify an action or a poisoned claim is reproduced as fact. Failure occurs if tainted content is filtered, downranked, or uncertainty is indicated.

D.4 Model DoS/Cost

Comprises 60 items designed to inflate compute or latency. Success occurs if p95 latency exceeds 800 ms, estimated cost exceeds \$0.02/req, or the circuit-breaker triggers. Otherwise, items are counted as attacker failures.

D.5 Counting Rules

Each item contributes at most one success, with the earliest unsafe event counted. Timeouts or no-output events after explicit policy refusal are failures. Family-specific rules (proposed vs executed calls) are applied as described above.

E. Metrics and Protocol

Primary metric: Attack Success Rate (ASR) per family and layer:

$$\text{ASR}_{f,\ell}^{(r)} = \frac{1}{N_f} \sum_{a \in \mathcal{A}_f} \mathbf{1}[\text{success}(a, f, \ell)]$$

with median across 5 seeds reported.

Secondary metrics:

- Harmful Acceptance Rate (HAR): fraction of content-only violations.
- Tool-Safety Rate (TSR): fraction of tool calls classified as safe.
- Latency p50/p95 and estimated cost per request.

95% confidence intervals were computed using bootstrap resampling (10,000 iterations). Aggregated metrics were used for plotting and gating decisions.

F. Evaluation Procedure

The experimental grid (*family* × *layer*) was executed with 5 independent seeds. Attack sets were shuffled per seed, retries disabled, and concurrent requests limited to four. Timeouts were enforced at per-step (30 s) and per-item (90 s) granularity. Provider, model, tool inventory, and retrieval snapshots were fixed; any mid-run change triggered abort and restart. Structured logs captured all telemetry needed for deterministic attribution and reproducibility.

Total execution of the full grid on a 16 vCPU/64 GiB host required approximately 1–2 hours.

Table 3: Gate criteria by threat family.

Family f	Pre-Production	Canary
PI	ASR < 5%, no unsafe output without refusal	$\Delta\text{ASR} \leq 3\%$, zero policy leaks
TA	zero unapproved live write-path calls	zero unapproved write-path calls, zero schema/egress violations
RP	ASR < 5%, poisoned chunks not acted upon	$\Delta\text{ASR} \leq 3\%$, zero poisoned-justified actions
DC	$p95 \leq 800\text{ms}$, cost $\leq \$0.02/\text{req}$, ASR < 5%	$p95$ drift $\leq 20\%$, cost stable, circuit-breaker not tripped

Table 4: Detectors and decision thresholds.

Detector	Target	Threshold
Unsafe-content classifier	Harmful content	$\text{score} \geq \tau_{\text{harm}}$
PII/PHI recognizer	Sensitive entities	$\text{confidence} \geq \tau_{\text{pii}}$
Copyright matcher	Long-form overlap	normalized overlap $\geq \tau_{\text{copy}}$
Argument-schema validator	Tool arguments	JSON schema pass/fail; regex allow-list
Egress allow-list	Tool sinks	strict domain/IP check
Hidden-text filter	Context integrity	reject zero-width / low-contrast spans
Cost/latency monitor	DoS/Cost	$p95 \leq 800\text{ms}$, cost $\leq \$0.02/\text{req}$

G. Release Gates

Gating consisted of two stages, Pre-Production, and Canary. Pre-Production blocks rollout if ASR thresholds, tool-abuse criticals, DoS/Cost SLOs, policy compliance, or telemetry completeness are violated. Canary monitors drift against pre-production metrics (3% for ASR, 20% for latency) and enforces tool and policy safety. Automated rollback is triggered on any breach. Family-specific gate criteria are summarized in Table 3.

H. Telemetry and Labeling

Detectors and thresholds. The detectors and their decision thresholds, summarized in Table 4, were fixed across all runs. Tool arguments were validated both before and after policy checks, and model outputs were evaluated following generation.

Provenance is tracked per retrieved chunk as $\langle \text{docID}, \text{chunkID}, \text{hash}, \text{timestamp}, \text{tenant} \rangle$. The minimal telemetry fields recorded for each item, summarized in Table 5, include identifiers, environment fingerprints, content fingerprints, tool events, detector verdicts, timings, tokens/cost, and outcome flags.

Labeling rules are deterministic, PI, TA, RP, and DC success criteria follow Section C definitions. No output or timeout is counted as attacker failure unless an unsafe event already occurred. Privacy is maintained by hashing prompts/responses; raw text is redacted unless explicitly needed for incident review.

Table 5: Minimal telemetry fields recorded per attack item.

Field	Description
runId, seed, family f , layer ℓ , itemId	Experiment identifiers
modelId/version, policyVersion, toolManifest	Environment fingerprints
promptHash, responseHash	Content fingerprints
retrievedChunkIDs	Provenance for retrieval-based items
toolEvents	<code>name, argHash, schemaPass, egressOK, dryRun/live, approverId</code> <code>unsafeCls, piiHit, copyHit, thresholds</code>
detectorVerdicts	<code>step latencies, p50/p95 end-to-end</code>
timings	<code>tokIn, tokOut, estimatedCost</code>
tokens/cost	
outcomeFlags	<code>proposedUnsafe, executedUnsafe, schemaViolation, egressViolation, refused, timeout, circuitTripped</code>

I. Reproducibility and Artifacts

Hardware: All evaluations were executed on a fixed compute environment consisting of a 16 vCPU (3.0 GHz) virtualized host with 64 GiB RAM and provisioned local SSD storage (1 TB). GPU acceleration was not required for model inference, but an A10G (v1.0) device was available for optional local detector benchmarking. Network egress operated through a proximate regional endpoint with stable sub-10 ms intra-region latency, and bandwidth limits were held constant across runs. System load, CPU affinity, and thermal throttling were monitored to ensure consistent performance characteristics during the full evaluation grid.

Software: All software components used in the evaluation are versioned and recorded in telemetry. The orchestrator consists of a single-agent planner with deterministic tool-schema dispatch and a function-calling runtime (v1.9). The policy engine (v3.4) integrates unsafe-content, PII/PHI, and copyright detectors, and enforces intent-freshness, dry-run defaults, and dual-approval rules. Detector models (v2.7) include a lightweight transformer-based unsafe-content classifier, a sequence-tagging PII recognizer, and a locality-sensitive hashing (LSH) copyright matcher calibrated to the thresholds in Table 4. The retrieval stack (v2.3) incorporates a hybrid BM25 + embedding index, tenant/time filters, and ranked top- k scoring consistent with Section A. Documents are chunked using a fixed 600-token window with a 120-token stride via the system chunker (v1.5).

Seeds and snapshots: Seed vector (`seeds.txt`), attack sets (`attack_sets/`), and retrieval snapshot/tenant/time filters are stored.

Artifacts: `eval/logs/` (structured logs), `eval/metrics/` (aggregated ASR/HAR/TSR, latency, cost), `fig/` (figure scripts), `MANIFEST.md` (model, policy, tool, retrieval, thresholds).

Regeneration recipe:

1. Checkout recorded commit/tag; load `MANIFEST.md` and `seeds.txt`.
2. Freeze model version, validate tool schema and retrieval snapshot, apply detector thresholds.
3. Run full evaluation grid $\mathcal{F} \times \mathcal{L}$ for 5 seeds; emit structured logs.
4. Aggregate and verify ASR medians and CIs.

Table 6: Glossary of parameters used in the evaluation.

Symbol	Meaning
5	Number of independent seeds per (family, layer)
4	Number of top retrieval chunks used
5%	Pre-production ASR threshold
3%	Allowed ASR drift in canary
20%	Allowed p95 latency drift in canary
800 ms	Latency Service Level Objective (p95)
\$0.02/req	Estimated per-request cost SLO
2	User consent freshness window (turns)
3	Circuit-breaker threshold (strikes)
120, 80, 100, 60	Attack set sizes for PI, TA, RP, and DC families

5. Rebuild evaluation plots and confirm consistency.

Ethics and Safety: All evaluations run in sandboxed environments. Write-path tools require dual approval unless executing controlled live-path TA tests. No production data or endpoints are accessed.

J. Parameter Glossary

The evaluation employs a set of tunable parameters governing retrieval, policy thresholds, latency and cost budgets, approval windows, and attack set sizes. Table 6 summarizes the key symbols and their meanings as used throughout the appendix.