



Nevriq Technologies

LLM Production Readiness Checklists

Enterprise RAG Systems & Agentic Workflow Automation

By Nevriq Technologies • <https://www.nevriq.com>

Free Download Edition • 2026-01-26

Common Core — Ship / No-Ship Gates (All LLM Systems)

NO-SHIP if any fail:

- AuthN/AuthZ in place for users and service-to-service
- PII and secret handling: redaction in logs, encrypted at rest and in transit
- Model and prompt versioning: pinned versions with rollback capability
- Observability: latency, error rate, cost, saturation dashboards and alerts
- Load and soak tests cover expected peak traffic
- Incident runbooks exist and are reachable by on-call teams

RAG Production Readiness Checklist

A) Data & Access Control (NO-SHIP)

- ACL-correct retrieval (doc-level and section-level where required)
- Source allowlist enforced (only approved repositories/indexes)
- Ingestion integrity validated (parsing, chunking, metadata)
- Freshness SLA defined and monitored
- Deletion propagation tested (index purges within defined window)
- Auditability: source_ids stored per response without raw sensitive text

B) Retrieval Quality (NO-SHIP)

- Golden set exists (real questions with expected sources)
- Retrieval recall@k meets threshold for critical categories
- Reranking evaluated with measurable improvement
- Context assembly rules defined (dedupe, diversity, recency bias)
- Query rewriting tested for regressions

C) Grounding & Answer Correctness (NO-SHIP)

- Citations required for factual claims
- Citation precision verified
- No-evidence behavior implemented (abstain / clarify / route)
- Hallucination rate below threshold on golden and adversarial sets
- Policy conflict handling defined with provenance rules

D) Prompt Injection Defense (NO-SHIP)

- Retrieved content treated as untrusted data
- Prompt-injection test suite passes
- Tool use disabled or tightly constrained for RAG-only systems
- Output filtering for data leakage patterns

E) Reliability & Degradation (NO-SHIP)

- Vector DB outage behavior implemented
- Low-confidence escalation path defined
- Safe caching policy documented
- Latency budgets enforced

F) Operations & Governance (NO-SHIP)

- Re-index cadence and backfill plan documented
- Monitoring includes retrieval latency, hit rate, zero-hit queries
- Human feedback loop implemented and triaged
- Red-team cadence established

Agentic Production Readiness Checklist

A) Tooling Safety & Permissions (NO-SHIP)

- Tool allowlist and explicit schemas enforced
- Least-privilege permissions per tool and action
- Two-phase commit for writes (plan → confirm → execute)
- Idempotency keys for all write actions
- Tool contract tests (timeouts, errors, schemas)
- Audit log for every action

B) Identity, Authorization & Abuse (NO-SHIP)

- Verified identity for sensitive actions
- Rate limiting and abuse detection
- Fraud/ATO heuristics with forced escalation
- PII minimization and masking

C) Agent Decision Quality (NO-SHIP)

- Intent routing accuracy meets threshold
- Entity extraction validated
- Tool selection correctness meets threshold
- Wrong-action rate below strict limits
- Clarifying-question behavior validated

D) Safety Policies & Escalation (NO-SHIP)

- Clear escalation triggers defined
- Safe refusal handling implemented
- User-visible receipts for all actions
- Explicit confirmation for irreversible actions

E) Reliability & Degradation (NO-SHIP)

- Tool outage degradation plan implemented
- Circuit breakers and safe retries
- Partial failure handling defined
- Latency budgets enforced
- Automatic safe-mode disables writes on anomaly

F) Observability & Postmortems (NO-SHIP)

- Distributed tracing across orchestrator and tools
- Action anomaly alerts
- Replayable traces with redaction
- Post-incident review process defined

G) Change Management (NO-SHIP)

- Prompt, model, and tool versions pinned
- Canary releases with rollback
- Sandbox/staging environment for tool execution
- Human approval for policy changes

Technical support: contact@nevriq.com

Subscribe to Nevriq Insights: <https://www.nevriq.com/solutions/support-copilot>

© Nevriq Technologies — Free distribution permitted with attribution.