

100 Controls For Agentic Al Security



Governance and Risk Management

Implementation: assign an executive owner; publish a charter; tie KPIs to risk reduction. Mitigates: unclear ownership, drift.

1. Executive accountability for Agentic Al

5. Policy set for Agentic Al

human-in-the-loop.

decisions.

outputs.

Implementation: publish policies for autonomy

limits, tool access, data usage, model use, and

Mitigates: inconsistent behaviors, misuse.

9. Safety-security trade-off review

Implementation: joint review of guardrails,

content filters, and mission goals; document

behavior); maintain examples. Mitigates: blind spots, inconsistent triage.

2. Agentic AI risk taxonomy

Implementation: define risk types specific to agents (prompt injection, rogue tools, model theft, emergent

Implementation: STRIDE+ATLAS profiling for planners, tools, memory, connectors; update each release. Mitigates: unaddressed attack paths.

3. System-level threat modeling for

agents

exceptions. exposure.

Mitigates: shadow go-lives, unmanaged

Implementation: formal sign-off for residual Agentic

Al risks with business justification and time-boxed

4. Risk acceptance workflow

8. Third-party and supplier assurance for Al **Implementation:** Al-specific TPRM questions for

6. Al use case intake and approval 7. Security requirements traceability Implementation: gated intake with risk scoring, **Implementation:** baseline control requirements PII sensitivity, tool reach, external calls, and blast per use case; track to tests and evidences. models, vector DBs, tool APIs, brokers; require pen test and SBOM. Mitigates: risky deployments, scope creep.

Mitigates: supply chain compromise. Mitigates: control gaps.

12. KPI and KRIs for AI risk

Implementation: define leading indicators like prompt-inj attempts blocked, unsafe tool calls averted, model exfil attempts. Mitigates: unknown risk posture.

Data Security and Privacy

Mitigates: over-permissive prompts, unsafe Mitigates: runaway actions, unsafe self-delegation.

radius.

Implementation: classify agent autonomy from 0 to 3; require stronger controls for higher tiers.

10. Autonomy tiering and approval

Implementation: add AI checkpoints into SDLC gates for data, prompts, tools, evals, and red teaming.

11. Secure development lifecycle for

Mitigates: late discovery of risks.

Agentic Al

Implementation: label inputs, memory, logs,

13. Data classification for AI flows

embeddings; block restricted classes from model context. Mitigates: accidental disclosure.

17. Embedding store protections

Implementation: encrypt at rest, row-level ACLs

Implementation: enforce least data in prompts and memory; token-budgeted masking.

14. Context minimization

Mitigates: leakage via responses or tool calls

18. No-train and retention controls

Implementation: set model flags or proxies to

disable training on production conversations;

Mitigates: data exposure, model memorization.

19. Confidential compute for sensitive

Implementation: TEEs or VPC-scoped inference

masking, or tokenization; test with seeded data.

15. PII and secret scrubbing

Implementation: pre-prompt scrubbing,

Implementation: curated datasets, poisoning detection, lineage tracking, and quarantine of untrusted data.t Mitigates: data poisoning, bias drift.

16. Training and fine-tune hygiene

agent runs; audit epsilon budgets.

by tenant and purpose, query-time filters. enforce retention windows. Mitigates: cross-tenant data bleed. Mitigates: unintended model learning.

21. Watermarking and content authenticity **Implementation:** attach provenance metadata to Al outputs; verify inbound content signatures

data sources agents can query. Mitigates: privilege creep.

Implementation: quarterly access review for all

22. Dataset access reviews

Mitigates: cloud insider and co-tenancy risk.

inference

for regulated workloads.

20. Differential privacy for analytics

Implementation: DP libraries for reporting on

26. Output handling guardrails

Implementation: schema validation, type

Mitigates: insecure output handling.

checks, regex safelists, and allow-only function

Mitigates: re-identification.

where supported.

Mitigates: spoofing, deepfake input.

Model and Prompt Security

rollback; map prompts to use cases. Mitigates: prompt drift, hidden backdoors.

Implementation: central store, approvals, diffs,

27. Model selection policy

31. Eval suite for security

23. Prompt registry with version control

guarantees, eval results; disallow unknown origins. Mitigates: model risk, licensing traps.

Implementation: continuous tests for injection, data

exfil, tool abuse, hallucinated APIs; publish scores.

Implementation: choose models by trust level, privacy

24. Prompt hardening patterns

Implementation: instruction isolation, deny-lists,

meta-prompt self-checks, tool preconditions.

Mitigates: prompt injection, jailbreaks.

28. Sensitive capability gating

Implementation: require explicit flags for code

execution, file write, or network calls; add human

approval for Tier 3 autonomy. Mitigates: unsafe actions.

32. Temperature and tool-use limits

Implementation: cap randomness and chain

Mitigates: instruction override.

25. Anti-prompt injection filters

Implementation: pre- and post-processing to

neutralize injection tokens, system override

attempts, and tool-call coercion.

29. Model redaction responses 30. Model access isolation **Implementation:** forced refusals for sensitive Implementation: per-use case API keys, separate queries; templated safe-completion fallbacks. projects, and VPCs; no shared tokens.

Mitigates: data leakage, misuse. Mitigates: blast radius.

names.

Mitigates: regression risk.

34. Tool permission contracts 33. Tool registry and attestation 35. Dry-run and simulation mode

function; enforce preconditions and max impact.

Implementation: granular scopes per tool

Mitigates: unknown capabilities.

Implementation: inventory every tool with

purpose, inputs, side effects, auth model, and

37. Network egress controls for tools **Implementation:** proxy allow-lists, DLP on egress, per-tool routing

41. Tool output taint tracking

Implementation: per-tool quotas, concurrency caps, cost ceilings, circuit breakers.

Mitigates: model DoS, bill shock.

38. Rate limiting and cost guards

42. Tool health and integrity probes

Implementation: heartbeat checks, checksum of

tool binaries or endpoints, version pinning.

Mitigates: supply chain tampering.

39. Command and code execution sandboxes

filesystem jails, outbound blocks by default.

Implementation: ephemeral containers, syscall filters,

Implementation: simulate tool calls with

Mitigates: unsafe execution.

Mitigates: RCE fallout.

to expected.

harmless stubs in test and lower envs; compare

Mitigates: irreversible changes.

36. Human-in-the-loop approval for

Implementation: explicit confirmation UI with

destructive actions

diff, impact estimate, and rollback plan.

40. Tool input validation and canonicalization **Implementation:** strict schema, path normalization, escaped shell args.

Mitigates: injection to downstream systems.

Mitigates: data exfiltration, C2.

risk rating.

Mitigates: chained injection.

Implementation: mark untrusted outputs; require sanitization before reuse in prompts or other tools.

sub-agents, task scope, and hand-off approvals.

Mitigates: privilege escalation via delegation.

only at start.

49. Constraint reminder meta-prompts

46. Memory governance

50. Safety critic agent

content.

Implementation: TTLs for episodic memory,

Mitigates: long-term leakage, bias lock-in.

quality scoring, and removal of toxic or sensitive

43. Planner constraints

44. Role-based multi-agent segregation 45. Delegation policy **Implementation:** explicit rules for spawning Implementation: cap steps, recursion, and Implementation: separate planner, critic,

executor identities with scoped capabilities.

Mitigates: single-agent compromise impact.

48. Reward shaping safeguards

Mitigates: misaligned optimization.

Implementation: penalize risky tool use, unsourced

claims, and privacy violations in feedback loops.

Mitigates: runaway plans, lateral risk.

47. Goal alignment checks

51. Explainability snapshots

map to stated objectives; abort on deviation. Mitigates: off-policy actions.

Implementation: pre- and mid-run checks that steps

parallel branches; enforce goal and scope limits.

Implementation: immediate halt command, revoke tokens, cancel queued tool calls.

52. Kill-switch and safe stop

Implementation: independent reviewer agent Implementation: inject constraints at each step, not with veto on sensitive outputs and actions. Mitigates: unsafe completions, tool misuse. Mitigates: erosion of instructions.

access

Implementation: short-lived tokens bound to a single run and tool; automatic rotation. **Mitigates:** token replay, sprawl.

execute code or handle sensitive data.

Mitigates: tampering, side-loading.

54. Ephemeral, scoped credentials

Identity, Secrets, and Access

Implementation: brokered secrets issuance per step; no secrets in prompts or memory. Mitigates: leakage and theft.

57. Secrets isolation and just-in-time

62. Security analytics for AI signals

Implementation: SIEM use cases for injection

Mitigates: stealthy attacks.

compliance; gate delivery.

Mitigates: harmful or illegal outputs.

patterns, unusual tool sequences, exfil indicators.

58. Device and environment attestation

Implementation: require verified runtime for agents that

Monitoring, Telemetry, and Detection

55. Delegation chain tracking

purpose across sub-agents and tools.

Mitigates: non-repudiation gaps.

approvals; signed attestations.

Mitigates: social engineering.

Mitigates: active exploitation.

Mitigates: supply chain tampering.

configs at runtime.

Implementation: verify tool endpoints, binaries, and

Implementation: propagate caller identity and

56. Policy-as-code for agent permissions

Implementation: OPA or ABAC policies referencing run

context, data class, and autonomy tier.

Mitigates: static RBAC over-grant.

61.Agent run logging **Implementation:** structured logs for prompts, tool calls, inputs, outputs, decisions, and approvals with hashes.

Mitigates: forensic gaps.

or atypical API regions.

67. Tool integrity monitoring 65. Cost and quota monitoring 66. Content risk scoring

Testing, Red Teaming, and Assurance

69. Security test plan per use case

70. Adversarial red teaming

Implementation: scheduled ATLAS-informed

Mitigates: real-world attack readiness.

74. Dataset health checks

Mitigates: training-time compromise.

artifacts; lock clean snapshots.

and model theft.

campaigns targeting injection, tool abuse, exfil,

Implementation: detect duplicates, outliers, poisoning

maintain blocklists and allowlists with metrics. Mitigates: brittle defenses.

Implementation: inject failures in model calls, tool

75. Pre-prod chaos drills

timeouts, and partial data.

Mitigates: resilience gaps

Incident Response and Resilience

80. Customer and regulator comms pack

76. Al-specific IR playbooks

Implementation: playbooks for injection, model

key theft, tool compromise, data leakage; include

policies, and tests; close the loop. Mitigates: recurrence.

automatically.

Implementation: feed lessons into prompts, tools,

77. Key rotation and model rebind

Implementation: rapid rotation for model and

tool creds; update config management

Mitigates: token replay after breach.

81.Post-incident eval and

guardrail tuning

82. Resilience testing of kill-switch **Implementation:** periodic drills to ensure safe-stop works under load.

78. Prompt and memory rollback

Implementation: versioned rollback of prompts.

guardrails, and memory stores.

Mitigates: fail-open risk.

Mitigates: persistent compromise.

Implementation: IaC and Git-Ops for prompts, tool policies, and guardrail configs.

Mitigates: release surprises.

Implementation: infrastructure-level quotas, per-project budgets, and pre-commit checks. Mitigates: economic DoS

services, brokers, and tools; rotate certs.

85. Secure inference perimeter

and controlled registries.

Mitigates: MITM.

Mitigates: model tampering.

Implementation: signed checkpoints, hashes,

89. TLS, mTLS, and cert pinning

Implementation: enforce mTLS between agent

91. Human oversight training

Implementation: teach reviewers how to spot injection, hallucinated tools, and unsafe plans. high-risk steps than those who authored prompts or tools. Mitigates: rubber-stamping risky actions. Mitigates: insider risk.

96. Accessibility and bias reviews

100. Continuous improvement council **Implementation:** monthly cross-functional review of metrics, incidents, red team learnings,

Implementation: keep artifacts for prompts, evals, red team results, incidents, and approvals.

97. Audit readiness pack

68. Provenance capture

and reasoning for critical outputs.

Mitigates: unverified decisions.

Implementation: store source links, evidence,

Supply Chain, Deployment, and Infrastructure

83. SBOM and MBOM for AI stack Implementation: IaC and Git-Ops for prompts, tool policies, and guardrail configs. Mitigates: drift, undocumented changes.

87. Environment parity

92. Separation of duties for approvals **Implementation:** different humans approve

Mitigates: failed audits, inability to prove control.

Implementation: run safety and security evals together Mitigates: safety holes created by security changes and vice versa.

79. Containment levers

Mitigates: active damage.

families, route to safe fallback models.

Implementation: throttle autonomy level, disable tool

86. Config and prompt as code

Mitigates: drift, undocumented changes.

90. Geo and residency controls

stores; block cross-region spillover.

Mitigates: regulatory breach.

Implementation: pin model regions and data

98. Decommissioning and model retirement

Implementation: revoke access, wipe memory

94. Usage transparency for end users

Implementation: disclose agent capabilities,

Mitigates: long-term leakage, bias lock-in.

data use, and escalation paths.

Copyright © Ampcus Cyber

84. Model artifact integrity **Implementation:** signed checkpoints, hashes, and controlled registries. Mitigates: model tampering.

Human Factors, Ops, and Compliance

88. Cost and quota guardrails in infra

93. Change management for prompts and tools **Implementation:** ticketing, peer review, and automated tests before merge Mitigates: accidental regressions.

> stores, archive artifacts, and update docs. Mitigates: orphaned risk.

and roadmap.

depth when tools are enabled Mitigates: erratic behavior, looping.

Tooling and Function Use Security

Mitigates: over-privilege.

Agent Orchestration and Autonomy

Implementation: capture plan graphs and decision rationales for high-risk runs. Mitigates: audit gaps, incident forensics friction. Mitigates: damage containment.

53. Agent identity distinct from users **Implementation:** separate service principals for agents; no reuse of human tokens. Mitigates: attribution gaps.

Mitigates: DoS, abuse, fraud.

Implementation: list of attack simulations,

coverage matrix for prompts, tools, and data

Implementation: alerts for cost anomalies, token floods,

73. Pen testing of tools and brokers Implementation: test function routers, sandboxes, and connectors like any other app surface. Mitigates: classic web and API flaws.

Mitigates: unknown weaknesses.

paths.

Implementation: pre-approved templates, artifact list, and evidence bundle for AI incidents. Mitigates: compliance and reputational impact.

roles and SLAs.

Mitigates: slow response.

Implementation: reproduce prod runtimes locally with masked data and stubbed tools.

95. Legal guardrails for output use Implementation: licensing checks, citation requirements, and high-risk domain disclaimers Mitigates: IP and regulatory violations.

99. User feedback loop and abuse

Implementation: in-product reporting, triage, and hotfix

reporting

pipelines for unsafe outputs.

Mitigates: prolonged exposure. Mitigates: stagnation and drift.

Implementation: test outputs for bias, fairness, and

accessibility for key user groups.

Mitigates: ethical and legal risk.

www.ampcuscyber.com

59. MFA proxies for human approvals 60. Outbound identity for external APIs Implementation: dedicated API identities and **Implementation:** step-up auth for destructive IPs; monitor for anomalous use. Mitigates: cross-system impersonation.

63. Model interaction anomaly detection 64. Data egress DLP for agents Implementation: baseline per use case; flag Implementation: pattern and policy checks spike in refusals, token usage, recursion depth. post-generation and pre-egress. Mitigates: data leakage.

72. Safety-security combined evals 71. Guardrail and filter testing Implementation: measure bypass rates; to catch trade-offs.