This side-by-side reference guide compares security controls for internal versus external MCP server deployments, with four priority ratings. The document includes shared concerns that apply to both models, quick validation questions for security assessments and a decision framework for choosing between self-hosted and vendor-hosted approaches. This guide is deal for security teams evaluating MCP architecture options or auditing existing implementations.

## Table 1: Priority Legend

| Priority | Level | Description |
|----------|-------|-------------|
| P1 | Severe | Must have before production. Fundamental security controls. |
| P2 | High | Should have for production. Important defense-in-depth. |
| P3 | Medium | Plan for implementation. Governance and operational maturity. |
| P4 | Recommended | Nice to have. Advanced capabilities for mature programs. |

## Shared Concerns

The concerns in Table 2 are essential regardless of whether you run internal or external MCP servers.

## Table 2: Shared Concerns

| Priority | Area | Control | Key Question |
|----------|------|---------|--------------|
| P1 | Kill switches | • Ability to quickly disable tools, endpoints or entire MCP connections | • Can you disable any agent/tool within five minutes? |
| P1 | Authentication | • Strong identity verification for all callers (humans, agents, services) | • Can an unauthenticated request ever reach the MCP? |
| P1 | Logging | • Capture prompts, tool calls and responses with appropriate redaction | • Can you reconstruct what an agent did and why? |
| P1 | Prompt injection defense | • Input sanitization<br>• Injection detection<br>• Output validation | • What happens if malicious content is in retrieved documents? |
| P2 | Human-in-the-loop | • Require approval for high-risk, irreversible or financial actions | • Which actions can an agent take without human approval? |
| P2 | Data classification | • Label and control data by sensitivity before it reaches MCP | • Do you know what data sensitivity levels the MCP can access? |

| Priority | Area | | |
|---|---|---|---|
| P2 | Tool minimization | • Only enable tools necessary for the use case; treat each tool as a risk | • Is there a tool enabled that isn't actively needed? |
| P3 | Lifecycle governance | • Catalog, risk-tier and manage MCP servers/connections through lifecycle | • Do you have a single source of truth for all MCP connections? |
| P3 | Incident response | • AI-specific playbooks for containment, evidence preservation and communication | • Do you have an incident response playbook specifically for AI/agent incidents? |
| P3 | Red-teaming | • Test for prompt injection, data exfiltration, tool abuse, jailbreaks | • When did you last red-team your MCP integrations? |

## Table 3: Internal vs. External MCP Servers

| Priority | Area | Internal MCP Server (Self-hosted) | External MCP Server (Vendor-hosted) |
|---|---|---|---|
| P1<br><br>P1 | Scope and risk | • Classify by data sensitivity, integration **tier** depth, agent capability (read-only → supervised → autonomous)<br>• Use 4-tier priority model | • Classify by what data you'll send and what actions the vendor can perform<br>• Start restrictive, expand only with justification |
| P1 | Network | • No public IPs<br>• Private subnets only<br>• Access via VPN/ExpressRoute<br>• Inbound only from gateway/bastion | • All traffic via controlled egress (API management/gateway)<br>• Allowlist vendor fully qualified domain names (FQDNs)<br>• Use private link if available |
| P1 | Identity and auth | • Entra ID + OAuth 2.1<br>• Validate JSON Web Token (JWT) issuer/audience exactly<br>• Short-lived tokens<br>• Token binding for replay prevention | • Prefer your identity provider (Entra) over vendor accounts<br>• Separate keys per environment<br>• Verify vendor token handling<br>• Confirm no sub-processor access |
| P1 | Authorization | • Enforce at gateway AND backend<br>• Object-level authz<br>• Map Entra claims to MCP roles (admin/owner/user/read-only) | • Map internal roles to vendor permissions<br>• Avoid broad admin scopes<br>• Require approval workflows for high-risk actions |
| P1 | Secrets and non-identities | • Unique service principal per **human** MCP/agent<br>• Managed identities preferred<br>• Vault for secrets<br>• Rotation + expiration alerts | • Keys only in vault, injected at gateway<br>• Never expose keys to users/agents.<br>• Rotate keys on schedule and staff departure |
| P1 | Input validation | • Direct and indirect prompt injection defense<br>• Input sanitization<br>• Output validation<br>• Schema validation<br>• Encoding checks | • DLP/redaction at gateway<br>• Prompt shielding<br>• Strip secrets before sending<br>• Context length limits |

| Priority | Category | Controls | Vendor considerations |
|---|---|---|---|
| **P2** | **Tool execution** | • Container/VM isolation<br>• Resource limits (CPU/memory/time)<br>• Network segmentation<br>• Filesystem isolation<br>• Recursive call limits | • Vendor tools are black boxes. Broker via your APIs.<br>• No direct write to core systems.<br>• Monitor for capability drift. |
| **P2** | **Data protection** | • Classify data<br>• Segment sensitive indexes<br>• Access checks at retrieval<br>• Data lineage<br>• Unlearning/deletion<br>• Poisoning detection | • Define allowed data classes<br>• Start with non-sensitive data<br>• Synthetic data for testing<br>• Minimal vendor logging<br>• Differential privacy |
| **P2** | **Session/memory** | • Encrypt context at rest<br>• Session timeouts<br>• Cross-session isolation<br>• Bound memory size<br>• Automatic clearing | • Verify vendor session isolation<br>• Confirm no cross-tenant leakage<br>• Test for context persistence across sessions |
| **P2** | **Gateway layer** | • Centralized gateway for all MCP access<br>• JWT validation at edge<br>• Rate limiting<br>• Schema validation<br>• Request signing | • Force all traffic through your gateway<br>• Inject vendor keys at gateway<br>• Log everything<br>• Rate limit per user/agent |
| **P2** | **Supply chain** | • Assess model dependencies, libraries, base images<br>• Monitor for vulnerabilities in tool integrations | • Map vendor's model providers and subprocessors<br>• Understand inference location<br>• Require notification of upstream changes |
| **P2** | **Logging and monitoring** | • Full tracing with redaction<br>• Log integrity (WORM)<br>• Oversight agents<br>• Metrics<br>• Correlation IDs<br>• Alert on anomalies | • Log everything on your side<br>• Ingest vendor signals<br>• Correlation IDs<br>• Behavioral baselines<br>• Response integrity monitoring |
| **P2** | **Availability/SLAs** | • Define RTO/RPO<br>• Backup configs and indices<br>• Geographic redundancy<br>• Disaster recovery testing | • Negotiate SLAs (uptime, latency)<br>• Degradation plans<br>• Failover options<br>• Monitor vendor status |
| **P3** | **Governance** | • Catalog MCP servers/agents<br>• Risk tier → capabilities<br>• Change management<br>• Promotion gates (shadow → autonomous) | • AI-flavored vendor risk assessment<br>• SOC 2/ISO<br>• Risk tier per vendor<br>• Concentration risk<br>• Financial health |
| **P3** | **Portability** | • Version control configs<br>• Documented deployment<br>• Reproducible infrastructure | • Assess lock-in<br>• API compatibility with alternatives<br>• Abstraction layers<br>• Migration runbooks |

| | | | |
|---|---|---|---|
| | | | • Data export |
| P3 | **Integrity verification** | • Request signing<br>• Replay prevention<br>• TLS everywhere, including internal | • Response authenticity<br>• Cert pinning<br>• Detect response modifications<br>• Monitor for vendor-side injection |
| P3 | **Legal/compliance** | • Data security posture management (DSPM) for data residency<br>• Regulatory mapping. Multi-tenancy controls if applicable. | • Data processing agreement (DPA)/business associates agreement (BAA)<br><br>• AI clauses (no training, intellectual property ownership)<br>• Cyber insurance<br>• Liability terms<br>• Regulatory fit |
| P3 | **Testing and validation** | • Threat model<br>• Red team<br>• Config review<br>• Continuous security testing in CI/CD<br>• Resource exhaustion tests | • Sandbox with synthetic data<br>• Red-team integration<br>• Continuous evaluation<br>• Canary deployments<br>• A/B testing<br>• Cross-tenant tests |
| P4 | **Multi-tenancy** | • Tenant isolation at MCP<br>• Scoped indices<br>• Per-tenant keys<br>• Per-tenant rate limits and audit logs | • Verify vendor tenant isolation<br>• Test for cross-tenant access<br>• Contractual isolation guarantees |

# Security Validation Questions

Use these questions to quickly assess MCP server security posture.

**Table 4: MCP Server Security Assessment Questions**

| Internal MCP Servers | External MCP Servers |
|---|---|
| Can requests reach MCP without going through the gateway? | Can any user/agent call the vendor directly (bypassing your gateway)? |
| What happens if a token is stolen? | Where does the vendor store your API keys and for how long? |
| Can one user access another user's data by manipulating IDs? | What data is the vendor allowed to log and retain? |
| What tools can an agent call without human approval? | Which of the vendor's tools can write to your systems? |

| | |
|---|---|
| How quickly can you disable a misbehaving agent? | How quickly can you cut all traffic to this vendor? |
| Where are secrets for this MCP server stored? | What happens to your data if the vendor is acquired? |
| Can you reconstruct what an agent did last week? | Does the vendor train models on your prompts/responses? |
| What's in your backup for this MCP and when was it tested? | What's your plan if this vendor has a major outage? |

## Table 5: When to Use Internal vs. External MCP Servers

| Favor Internal MCP Servers When… | External MCP Servers May be Acceptable When… |
|---|---|
| Processing PII, PHI, PCI or highly confidential data | Working with public data or low-sensitivity internal data |
| Regulatory requirements mandate data residency control | Vendor meets all compliance requirements with attestation |
| Need full control over model behavior and guardrails | Vendor guardrails are sufficient for use case |
| Actions have high financial or operational impact | Actions are read-only or easily reversible |
| Require deep integration with internal systems | Integration is limited to specific, well-bounded use cases |
| Long-term strategic capability requiring investment | Rapid experimentation or time-to-value is critical |