

Hybrid-Cloud AI Reference Architecture Blueprint

A Multi-Layer Enterprise Architecture for Agentic AI, MCP Integration,
and Hybrid Deployment

Digital Enterprise Architecture Advisors (DEAA)

Executive Summary

This blueprint defines a **complete enterprise reference architecture** for deploying AI systems across hybrid-cloud environments. It integrates IBM's agentic-AI infrastructure vision, the Model Context Protocol (MCP), modern data engineering patterns, and the Agent Development Lifecycle (ADLC). The architecture spans six layers — hybrid infrastructure, AI control plane, data and feature layer, model and inference layer, agentic AI layer, and cross-cutting observability and governance — providing a structural foundation for secure, governed, scalable enterprise AI.

Author

Christian Kobsa

Strategic Enterprise Architect
Digital Enterprise Architecture Advisors

Version

v1.0 — 2026

Table of Contents

Executive Summary.....	1
1. Technical architecture blueprint	3
1.1 Strategic context and executive summary	3
1.2 From applications to agents (agentic enterprise architecture)	4
1.3 Why hybrid cloud is the natural home for agentic AI	4
2. Core architectural principles for enterprise AI.....	5
2.1 Acceptable agency.....	5
2.2 Interoperability and open standards	6
2.3 Secure-by-design	6
2.4 Evaluation-first development	7
2.5 Hybrid deployment.....	7
2.6 Continuous governance.....	7
3. Enterprise AI reference architecture (hybrid cloud).....	8
3.1 Layer 1 — Hybrid cloud infrastructure fabric.....	8
3.2 Layer 2 — AI control plane (LLM gateway + MCP)	9
3.3 Layer 3 — Enterprise data & feature layer	10
3.4 Layer 4 — AI models & inference layer	11
3.5 Layer 5 — Agentic AI layer	11
3.6 Layer 6 — Observability, security & governance	12
4. End-to-end reference architecture (textual view).....	13

1. Technical architecture blueprint

Architecting a Corporate AI-Based Enterprise Architecture in a Hybrid-Cloud Environment

1.1 Strategic context and executive summary

Strategic shift:

Enterprises are moving from application-centric to agentic architectures where AI agents become the primary execution and decision layer across IT, business operations, and customer workflows.

Key drivers:

- Agentic AI:

AI agents that can perceive context, reason over goals, and act through tools with policy-aware autonomy.

- Model Context Protocol (MCP):

A standardized, secure interface for tool access and enterprise integration.

- Hybrid cloud:

Need to unify on-premises, private cloud, and public cloud AI services.

- Data engineering evolution:

Data engineering becomes a strategic discipline powering AI readiness, not just pipelines.

- DevSecOps-driven governance:

Continuous evaluation, security, and risk management for AI agents via an extended Agent Development Lifecycle (ADLC).

Goal of the blueprint:

Provide a complete enterprise reference architecture for a corporate AI ecosystem that is:

- Secure-by-design

- Governed end-to-end
- Hybrid-cloud native
- Agent-ready
- MCP-integrated

1.2 From applications to agents (agentic enterprise architecture)

IBM's view of AI agents:

AI agents are not standalone tools; they operate inside complex hybrid ecosystems and must deeply integrate with enterprise systems. Agents:

- Perceive context
- Reason over goals and constraints
- Act through tools and services
- Operate under policy-aware autonomy
- Require continuous evaluation and governance

Agent Development Lifecycle (ADLC):

Extends DevSecOps with:

- Behavioral evaluation (beyond unit tests)
- Guardrails (acceptable agency, safety thresholds)
- Observability of reasoning traces (agent transcripts, tool calls)
- Runtime optimization loops (quality and cost)
- Governance and certification (catalogs, approvals, risk posture)

1.3 Why hybrid cloud is the natural home for agentic AI

Hybrid cloud is essential because:

- Enterprise data is distributed across on-prem, private cloud, and SaaS.

- Regulatory and sovereignty constraints require local inference.
- Latency-sensitive workloads (manufacturing, finance, operations) need edge/on-prem execution.
- Cloud AI services provide elasticity and access to frontier models.

IBM's hybrid infrastructure footprint:

- IBM Cloud
- IBM Z
- IBM Power
- IBM Storage
- IBM TLS

These are becoming agentic AI-ready, with MCP servers exposing them as secure, standardized tools across hybrid environments.

2. Core architectural principles for enterprise AI

Across IBM, Databricks, Snowflake, and MIT research, six principles consistently emerge.

2.1 Acceptable agency

Agents must have bounded autonomy:

- Explicit authority scopes (what each agent is allowed to do)
- Human-in-the-loop escalation paths
- Kill switches for emergency shutdown
- Reversible actions where possible
- Immutable audit trails for all decisions and actions

2.2 Interoperability and open standards

MCP as backbone:

- Tool access (APIs, systems, services)
- Resource exposure (data, knowledge, operations)
- Prompt schemas (typed, structured)
- Cross-platform orchestration (multi-agent, multi-system)

Agents interact with:

- Cloud APIs
- On-prem systems
- Databases
- IT operations tools
- Business applications

through MCP servers, not ad-hoc integrations.

2.3 Secure-by-design

Security is embedded at every layer:

- Identity propagation for agents (agents have identities; actions are attributable)
- Sandboxing of tools and code execution (gVisor, Firecracker, seccomp, container profiles)
- Network isolation (zero-trust, segmentation)
- Policy enforcement at the gateway (AI control plane)
- Continuous red teaming (prompt injection, adversarial testing)
- Memory poisoning defenses (protect agent memory from malicious data)

2.4 Evaluation-first development

Agents require evaluation-first practices:

- Behavioral benchmarks (task success, trajectories)
- LLM-as-a-Judge scoring (LLM-aaJ)
- Drift detection (behavioral and performance drift)
- Hallucination metrics (context vs. output comparison)
- Bias and fairness checks (across groups and metrics)

2.5 Hybrid deployment

Agents must run:

- In cloud-native environments
- On-premises (IBM Z, Power, Storage, VMware, Kubernetes)
- At the edge (local inference nodes)
- Across multi-cloud providers

with consistent governance, identity, and observability.

2.6 Continuous governance

Governance spans:

- Data lineage
- Model provenance
- Tool catalogs
- Agent registries
- Compliance evidence
- Audit logs

and is continuous, not periodic.

3. Enterprise AI reference architecture (hybrid cloud)

The architecture is layered; each layer is logically distinct but tightly integrated.

3.1 Layer 1 — Hybrid cloud infrastructure fabric

Components:

- Public cloud:

IBM Cloud, AWS, Azure, GCP

- Private cloud:

OpenShift, VMware, IBM PowerVS

- On-prem systems:

IBM Z, storage arrays, databases

- Edge compute:

Edge nodes, gateways

Key capabilities:

- Secure connectivity:

VPN, Direct Connect, SD-WAN

- Identity federation:

OIDC, IAM, LDAP

- Network segmentation:

Micro-segmentation, zones

- Zero-trust access:

Least privilege, continuous verification

IBM enhancements (agentic AI-ready infrastructure):

IBM infrastructure is becoming agentic AI-ready, with MCP servers for:

- IBM Cloud
- IBM Storage Insights
- IBM PowerVS
- IBM Z (e.g., watsonx Assistant for Z)
- IBM TLS

These MCP servers provide standardized, secure access for agents across hybrid environments.

3.2 Layer 2 — AI control plane (LLM gateway + MCP)

This is the central nervous system of the enterprise AI architecture.

Functions:

- Policy enforcement
- Routing to models (cloud, on-prem, edge)
- Identity propagation (user → agent → tools)
- Logging and audit
- Rate limiting and throttling
- Guardrail enforcement (safety, content filters)
- Cost governance (budgets, usage tracking)

Why MCP matters:

- MCP servers expose enterprise systems as typed, governed tools.
- Agents interact through a single, secure interface, not direct API calls.
- Tool schemas are standardized and versioned.

Benefits:

-
- Eliminates credential sprawl
 - Standardizes tool schemas
 - Enables multi-agent orchestration
 - Provides auditability across all tool interactions

3.3 Layer 3 — Enterprise data & feature layer

Integrates insights from Snowflake, Databricks, and MIT research.

Capabilities:

- Unified data lakehouse
- Vector indexes for RAG (retrieval-augmented generation)
- Real-time pipelines (streaming, event-driven)
- ML feature stores
- Data governance (lineage, access control, quality)

Trends (from MIT / Snowflake research):

- Data engineers now spend 37% of their time on AI (up from 19%).
- Unstructured data and real-time ingestion are becoming dominant.
- AI-powered data engineering tools significantly accelerate productivity.

Implications for architecture:

The data layer must support:

- Multimodal data (text, images, logs, time series, etc.)
- High-throughput ingestion
- Streaming pipelines
- Automated quality checks
- Synthetic data governance

3.4 Layer 4 — AI models & inference layer

Hybrid model strategy:

- Cloud frontier models:

Claude, GPT, Gemini, watsonx, etc.

- On-prem models:

Llama, Granite, Mistral, others

- Domain-specific models:

Fine-tuned for industry or enterprise tasks

- Multimodal models:

Text, image, audio, video, structured data

Deployment options:

- Cloud inference via private endpoints (no public internet exposure)
- On-prem inference on IBM Z, Power, or GPU clusters
- Edge inference for low-latency use cases

Runtime optimization:

- Model routing (choose best model per task)
- Cost-aware inference (optimize for price/performance)
- Caching (responses, embeddings)
- Distillation and quantization (optimize smaller models)

3.5 Layer 5 — Agentic AI layer

Implements the Agent Development Lifecycle (ADLC).

Agent capabilities:

- Memory:

Short-term, long-term, episodic

- Planning:

ReAct, Reflexion, hierarchical planning

- Tool use via MCP:

Typed tools, governed access

- Multi-agent collaboration:

Delegation, coordination

- Self-evaluation:

Internal checks, LLM-aal, feedback loops

Agent types:

- IT operations agents (self-healing, incident response)
- Security agents (threat detection, policy enforcement)
- Customer service agents (support, case handling)
- Research agents (knowledge synthesis, analysis)
- Workflow automation agents (process orchestration)

Agent governance:

- Versioning (prompts, policies, tools)
- Risk posture classification (low/medium/high risk)
- Behavioral guardrails (acceptable agency)
- Certification before deployment (tests, red teaming, approvals)

3.6 Layer 6 — Observability, security & governance

Observability:

- Reasoning traces (agent thought processes)
- Tool call logs (inputs, outputs, timing)

- Latency and cost metrics
- Drift detection (behavioral and performance)

Security:

- Sandboxing:

gVisor, Firecracker, seccomp, container security profiles

- Network isolation:

Restricted egress/ingress, zero-trust

- Memory poisoning detection
- Prompt injection defenses

Governance:

- NIST AI RMF alignment
- Audit trails
- Compliance evidence
- Kill switches (per agent, per tool, per environment)

4. End-to-end reference architecture (textual view)

A simplified vertical flow:

1. Hybrid Cloud Infrastructure Fabric

→ provides compute, storage, networking, identity, and connectivity across public cloud, private cloud, on-prem, and edge.

2. AI Control Plane (LLM Gateway + MCP)

→ centralizes policy, routing, identity, guardrails, logging, and cost governance; connects to MCP servers.

3. Enterprise Data & Feature Layer

→ lakehouse, vector indexes, pipelines, feature stores, and governance for multimodal, real-time, and high-quality data.

4. AI Models & Inference Layer

→ hybrid portfolio of cloud, on-prem, and edge models, optimized via routing, caching, distillation, and quantization.

5. Agentic AI Layer

→ agents with memory, planning, tool use via MCP, multi-agent collaboration, and self-evaluation, governed via ADLC.

6. Observability, Security & Governance

→ cross-cutting layer providing traces, logs, metrics, sandboxing, isolation, NIST AI RMF alignment, audit trails, and kill switches