



😊 Add icon 🖼️ Add cover

Omni Partner Platform

▼ Document Control

Captures version history, ownership, approvers, and distribution so changes are auditable.

	≡ Field	≡ Value
1	Project	JioOmni
2	Doc Owner	 Ravindra Chhabra  Mujtaba Rasool
3	Version	0.1-DRAFT
4	Last Updated	
5	Approvers	
6	Distribution	

1. Overview

The Partner Platform is a unified ecosystem that enables external developers, enterprises, and technology partners to seamlessly integrate with Omni. It provides the tools, workflows, and governance required to create, test, and publish AI-powered **Skills** that can be consumed by Omni users.

Requirements

2.1 Core Objectives

- **Grow Developer Ecosystem** – Provide self-serve tools for building, testing, and deploying Skills.
- **Secure by Design** – Enterprise login, CAPTCHA, 2FA, and governed approval workflows.
- **Seamless Onboarding** – Allow users to join their existing teams or create new ones.
- **Quality Assurance** – Automated validation of MCP servers/OpenAPI/Agent specs and pre-publish testing.
- **Scalable Architecture** – Modular foundation supporting phased rollout and future Agent-to-Agent integration.
- **Collaborative & Transparent** – Dashboards, analytics, and feedback loops for partners and governance teams.

2.2 Non-Functional Requirements (NFRs)

Defines service-quality targets—latency, uptime, DR, accessibility—that shape infra sizing and SLAs.

	≡ Category	≡ Target
1	Latency (P95)	
2	Uptime (SLA)	
3	Throughput	

4	DR RPO/RTO	
5	Accessibility	

2.3 KPIs

Establishes measurable success metrics (e.g., skill-success %, CSAT) used for post-launch health checks.

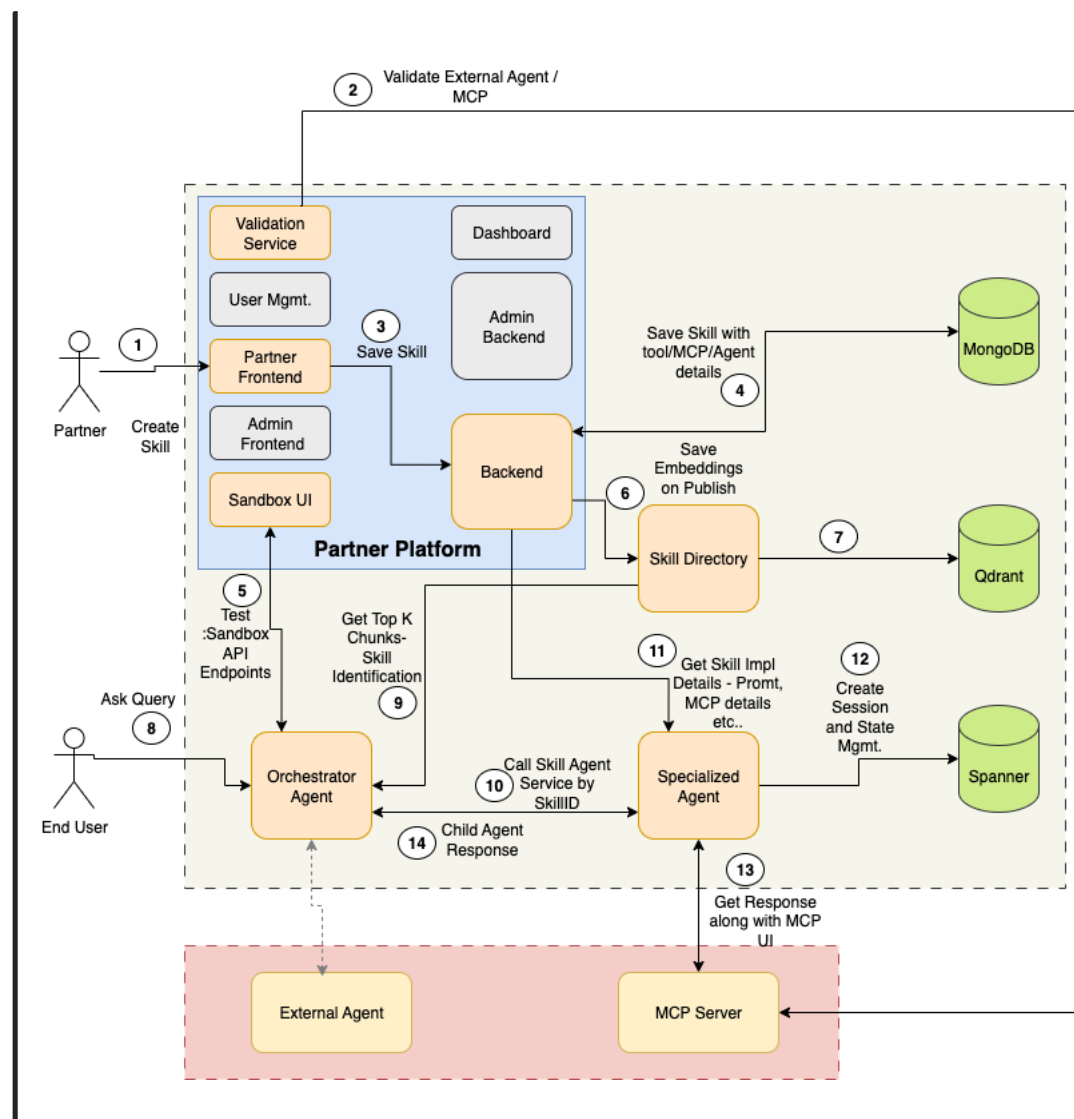
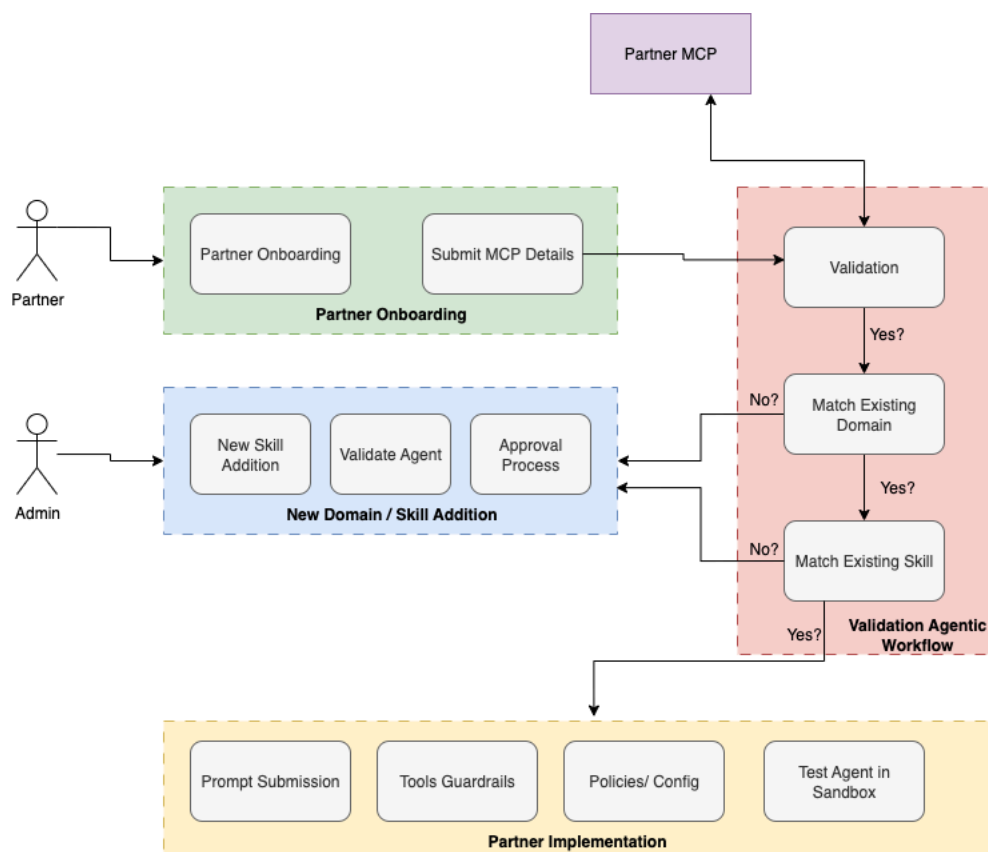
System Architecture

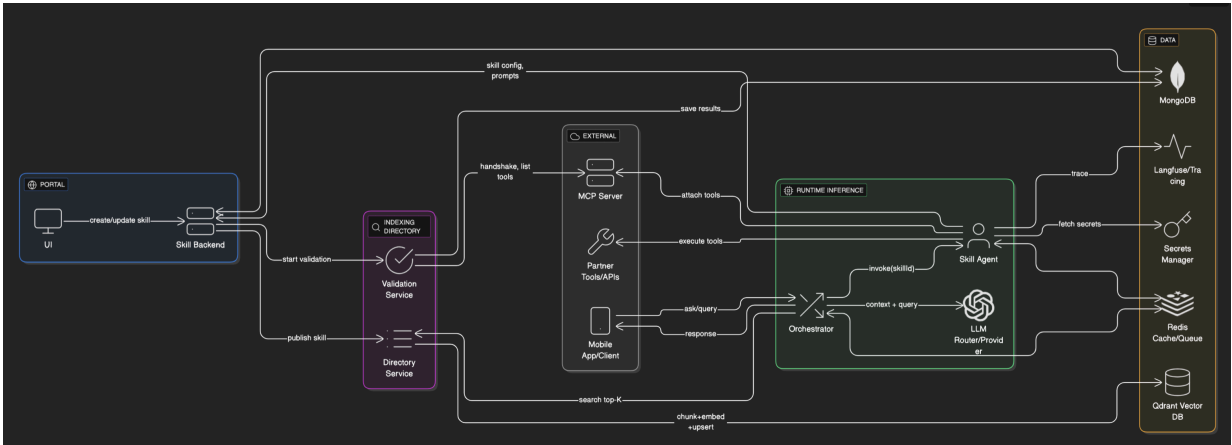
3.0 Diagram Index

Serves as a table of embedded Loop images so reviewers can find every architecture visual quickly. *(Just put which is important to your team)*

	≡ #	≡ Diagram	≡ Loop Embed / Link
1	1	High-Level Overview	
2	2	Deployment Topology	
3	3	Network Zoning	
4	4	Multi-Modal Flow	
5	5	Sequence Diagram	
6	

3.1 High-Level Overview (Diagram link)





3.2 Component Catalog

Summarizes each subsystem’s role and how it fits into the overall architecture.

	≡ Area	≡ Component	≡ Summary
1	Core	Supervisor/Orchestrator Service	Front-door for user queries (SSE/Socket.IO/HTTP). Retrieves top-K skills from Skill Directory, performs LLM routing, invokes Skill Agent Service, streams responses. Enforces retries/fallbacks and budgets.
2	Core	Skill Agent Service	Stateless runtime that materializes an Agent per request using skill prompts + binding. Attaches MCP tools, runs plan→act loop, calls external tools, streams results back to Orchestrator.
3	Core	LLM Router/Provider	Chooses best skill from retrieved candidates using query + snippets + policies; supports model routing and temperature/budget controls.
4	Core	Memory/Session Layer	Conversation/session context store (Spanner or Redis) with user context and tool results cache; optional long-term memory/personalization store.
5	Registries	Skill Directory Service (SDS)	Unified indexing & search: chunks skill text, generates embeddings, upserts vectors to Qdrant keyed by skillId/versionId; exposes /index & /search; applies tenant/visibility filters and health-based boosting.

6	Registries	Tools/MCP Registry	Catalog of internal MCP endpoints (baseUrl, transport, protocolVersion, capabilities, health) discovered via Validation Service; referenced by bindings.
7	Partner Portal	Portal Backend	Canonical skill metadata in MongoDB (prompts, versions, policies, bindings list); read API for runtime.
8	Validation & Indexing	MCP Validation Service	Performs handshake/auth checks, lists tools, captures schemas (and optional <code>x-*</code> policy hints), runs smoke tests, writes <code>validation_runs</code> + health signals.
9	Data & Storage	MongoDB	Source of truth for skills, versions, MCP endpoints, bindings, validation runs, audit trails.
10	Data & Storage	Qdrant Vector DB	<code>skills_v1</code> collection for retrieval; metadata filters on tenant/visibility/version/state; HNSW tuned for recall/latency.
11	Data & Storage	Google Secrets Manager/ Redis	Stores all partner credentials; app holds only <code>secret_ref</code> . Short-lived tokens fetched at runtime.
12	Data & Storage	Blob Storage (GCS/Azure Blob)	Stores large OpenAPI specs/manifests; referenced by <code>specRefs</code> .
13	Security	Policy/Guardrails	Prompt linting, tool allowlists, PII redaction; enforced in SAS and at response boundary.
14	Observability	Langfuse	Tracing of LLM/tool events with <code>conversation_id/request_id</code> ; links validation and runtime spans.
15	Observability	Log	System metrics/logs/traces to Databricks SLOs & alerts (p95 search, MCP health, publish pipeline).
16	Messaging & Events	Redis (Cache/Streams)	Low-latency cache for runtime skill reads and SDS candidate cache; ephemeral queues and Socket.IO adapter.
17	Client & UX	Partner Portal (Web UI)	Skill creation/edit, MCP registration, validation runs, publish/rollback, versioning, RBAC; shows discovered tools & health.

18	Client & UX	Client Web SDK (Sandbox)	Sends ASK events, receives streamed responses; carries user/session context and consent flags.
19	Governance	Audit & Approvals	Immutable logs of changes (who/what/when), approval workflow for publish, policy exceptions with expiry.

2.1 Partner Web App (Portal UI)

- Create/edit **Skill** (name, tags, domain, visibility, owners).
- Define **prompts** (system/instructions, tool-use guidelines, guardrails), **descriptions**, **test cases**.
- Register **MCP server(s)** (base URL, auth type, headers/keys referenced from Secrets Manager, not raw secrets in browser).
- Trigger **Validate** → displays discovered tools and validation status with logs.
- **Publish** workflow with version notes and rollout (draft → published → deprecated).
- Role-based access control (Partner Owner, Editor, Viewer).

2.2 Portal Backend Service (PBE)

- REST API for Skills, MCP endpoints, versions, validation runs.
- Persist canonical skill models in MongoDB.
- Starts validation jobs (enqueue) and embedding/index jobs.
- On **publish**, calls to **Skill Directory Service**.
- Exposes read API for **Skill Agent Service** (runtime fetch of skill details).

2.3 MCP Validation Service

- Performs connectivity handshake to partner MCP endpoints.
- Discovers tools, pulls tool schemas, runs **smoke tests** (optional sample tool calls with partner-provided test inputs).
- Caches tool list & availability, stores **validation reports** (status, errors, latency stats) in MongoDB.
- Periodic re-validation (cron/Temporal) & health signals for routing.

2.4 Skill Directory Service (SDS — Index & Search)

- Single service responsible for **chunking text**, **generating embeddings**, and **upserting** vectors to Qdrant against `skillId/versionId`.
- Exposes `POST /index` for publish/reindex (idempotent), and `POST /search` for runtime retrieval.
- Ranking (hybrid optional): cosine score + availability/health boost (+ optional BM25/text filters).
- Enforces tenant/visibility filters and policy gating (only `published`).
- **Index Map (Mongo)**: maintains mapping of `skillId + versionId → pointIds` for deterministic deletes/rollbacks and audit.

Indexes: `(tenantId, skillId, versionId) unique; active; batchId`.

VisualBasic

```
1 //Sample--
2 // collection: sds_index_map
3 {
4   "_id": "ObjectId",
5   "tenantId": "tenant-jio",
6   "skillId": "skill-astro-...",
7   "versionId": "sv-astro-1-0-0", // or semver "1.0.0"
8   "collection": "skills_v1",
9   "pointIds": [1200101, 1200102, 1200103], // Qdrant point IDs (int or UUID—pick one
10  and standardize)
11   "numPoints": 3,
```

```
11 "model": "text-embedding-3-large",
12 "chunking": { "algo": "recursive", "maxTokens": 512, "overlap": 64 },
13 "checksum": { "textHash": "sha256:...", "configHash": "sha256:..." },
14 "batchId": "idx-2025-09-04-astro-0001",
15 "active": true, // set false on rollback/delete
16 "createdAt": "2025-09-04T10:30:00Z",
17 "updatedAt": "2025-09-04T10:31:00Z",
18 "indexedBy": "svc-sds"
19 }
20
21
22
```

- Why this mapping?
- Guarantees **precise deletes** even if filters or metadata evolve.
- Enables **rollback** to a previous version by knowing exactly which point IDs belong to each version.
- Supports **auditing** (how many vectors, when, which model/chunking).

Clean delete/rollback options

- By IDs (preferred when map exists): `qdrant.points.delete(ids = pointIds)`.
- By filter (when no map): filter on `skill_id` + `version_id` payload keys.

API Contracts (concise) : **WS Wagisha Singh** can add for Skill Directory

2.6 Orchestrator Service (Supervisor)

- Receives `ASK` from clients (Socket.IO/SSE/HTTP), manages conversation context.
- Calls **Skill Directory Service** for top-K candidate skills.
- LLM-based **routing**: feeds retrieved skill snippets + user query + policy constraints to decide target skill.
- Calls **Skill Agent Service** with `skill_id` (and optionally version pinning), passes condensed query & user context (optional).
- Streams partial/full responses to client; handles retries/fallbacks.

2.7 Skill Agent Service (SAS)

- Stateless runtime that **materializes an Agent on demand**:
 - Fetches skill prompts/config from PBE- Portal backend.
 - Fetches secrets (e.g., auth tokens) from Google Secret manager by secret refs only.
 - Connects to MCP server(s), registers exposed tools as function/tool adapters.
 - Runs the agent loop (planning → tool calls → synthesis), with timeouts, rate limits, guardrails.
 - Streams tokens/results back to Orchestrator.
- Traces all LLM/tool events to **Langfuse**.
- Enforces per-partner quotas, concurrency, and per-tool policy.

2.8 Data Layer

- **MongoDB:** canonical skill data, versions, MCP endpoints, validation runs, publication history, policy flags, audit logs.
- **Qdrant:** vector embeddings for retrieval; collection `skills_v1` with metadata for tenant/skill/version/state/tags.
- **VAULT:** all secrets (MCP auth, API keys). Backend stores only references.
- **Redis:** low-latency caches (skill read cache), streaming buffers, and queues if needed.

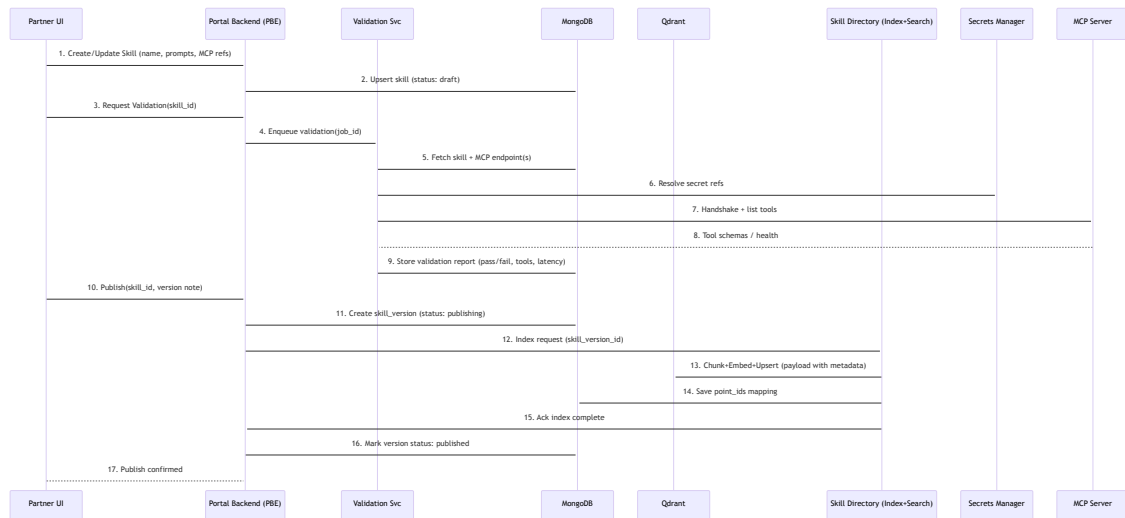
2.9 Observability & Safety

- **Langfuse** for tracing of LLM & tool calls (link with conversation_id, request_id).
- Safety/guardrails: prompt linting, allow/deny tool list per skill, response redaction.

3) Key Sequences

3.1 Partner: Create → Validate → Publish Skill

```
1 sequenceDiagram
2 participant U as Partner UI
3 participant B as Portal Backend (PBE)
4 participant V as Validation Svc
5 participant DB as MongoDB
6 participant Q as Qdrant
7 participant S as Skill Directory (Index+Search)
8 participant Va as Secrets Manager
9 participant M as MCP Server
10
11 U->>B: 1. Create/Update Skill (name, prompts, MCP refs)
12 B->>DB: 2. Upsert skill (status: draft)
13 U->>B: 3. Request Validation(skill_id)
14 B->>V: 4. Enqueue validation(job_id)
15 V->>DB: 5. Fetch skill + MCP endpoint(s)
16 V->>Va: 6. Resolve secret refs
17 V->>M: 7. Handshake + list tools
18 M-->>V: 8. Tool schemas / health
19 V->>DB: 9. Store validation report (pass/fail, tools, latency)
20 U->>B: 10. Publish(skill_id, version note)
21 B->>DB: 11. Create skill_version (status: publishing)
22 B->>S: 12. Index request (skill_version_id)
23 S->>Q: 13. Chunk+Embed+Upsert (payload with metadata)
24 S->>DB: 14. Save point_ids mapping
25 S->>B: 15. Ack index complete
26 B->>DB: 16. Mark version status: published
27 U-->>B: 17. Publish confirmed
28
29
```

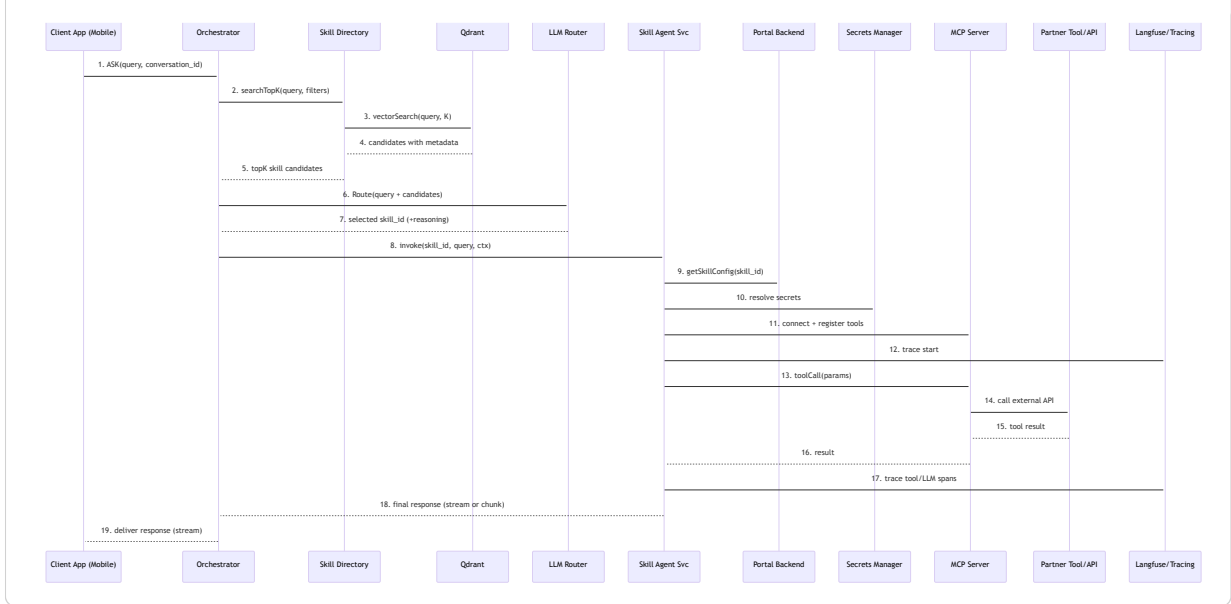


3.2 End-User: Query → Route → Invoke Skill → Tools → Answer

```
1 sequenceDiagram
2 participant C as Client App (Mobile)
3 participant O as Orchestrator
4 participant S as Skill Directory
5 participant Q as Qdrant
6 participant L as LLM Router
7 participant A as Skill Agent Svc
8 participant B as Portal Backend
9 participant V as Secrets Manager
10 participant M as MCP Server
11 participant T as Partner Tool/API
12 participant F as Langfuse/Tracing
```



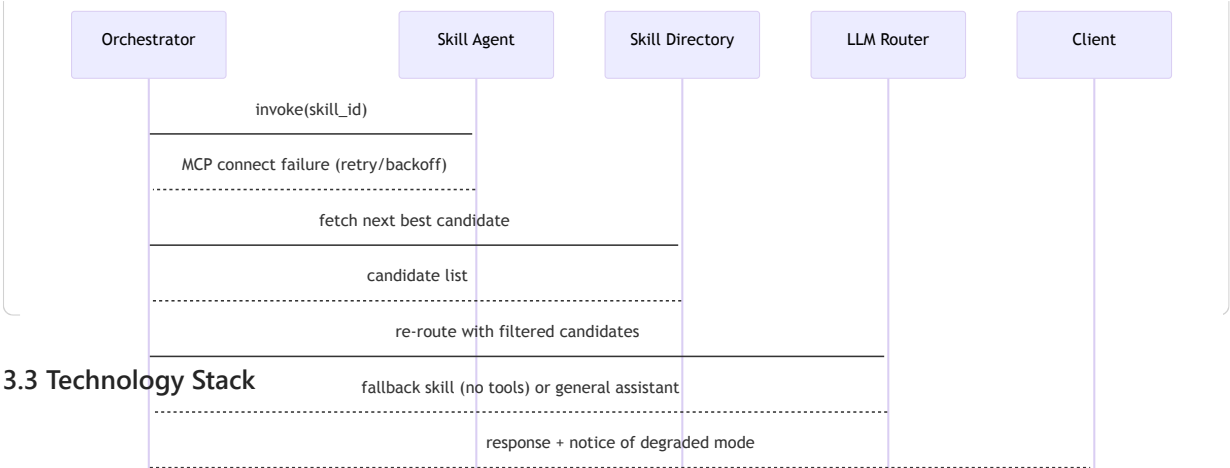
```
13
14 C->>O: 1. ASK(query, conversation_id)
15 O->>S: 2. searchTopK(query, filters)
16 S->>Q: 3. vectorSearch(query, K)
17 Q-->>S: 4. candidates with metadata
18 S-->>O: 5. topK skill candidates
19 O->>L: 6. Route(query + candidates)
20 L-->>O: 7. selected skill_id (+reasoning)
21 O->>A: 8. invoke(skill_id, query, ctx)
22 A->>B: 9. getSkillConfig(skill_id)
23 A->>V: 10. resolve secrets
24 A->>M: 11. connect + register tools
25 A->>F: 12. trace start
26 A->>M: 13. toolCall(params)
27 M->>T: 14. call external API
28 T-->>M: 15. tool result
29 M-->>A: 16. result
30 A->>F: 17. trace tool/LLM spans
31 A-->>O: 18. final response (stream or chunk)
32 O-->>C: 19. deliver response (stream)
33
34
35
```



3.3 [WIP] Failure: MCP Unavailable → Graceful Degradation (To be Explored more)

</> Mermaid

```
1 sequenceDiagram
2 participant O as Orchestrator
3 participant A as Skill Agent
4 participant S as Skill Directory
5 participant L as LLM Router
6
7 O->>A: invoke(skill_id)
8 A-->>O: MCP connect failure (retry/backoff)
9 O->>S: fetch next best candidate
10 S-->>O: candidate list
11 O->>L: re-route with filtered candidates
12 L-->>O: fallback skill (no tools) or general assistant
13 O-->>Client: response + notice of degraded mode
14
15
```



3.3 Technology Stack

	≡ Layer	≡ Technology	≡ Rationale
1	Validation Service	Go/Java	Efficient concurrent network probes (handshake, list_tools, smoke), small binaries, predictable latency.
2	Cache / Realtime	Redis (cache + pub/sub/streams)	Low-latency bootstrap caches, SDS candidate cache, SSE/ WebSocket adapters .
3	Messaging / Events	GCP Pub/Sub	Decouples validation/indexing/publish pipelines; durable fan-out and replay for audits.
4	Secrets	Google Secret Manager	Store only <code>secret_ref</code> (no plaintext); rotation, audit trails.
5	LLM Provider Interface	Vendor-agnostic HTTP adapter (LiteLLM) - Router	Swap models/providers behind a stable interface; record tokens/latency.
6	Observability — Tracing	Langfuse	End-to-end traces for LLM/tool spans; correlate validation ↔ runtime.
7	Observability — Metrics/Logs	Google Logs	p95 search latency, MCP health, publish pipeline SLOs; centralized logs.
8	Frontend — Partner Portal	React + TypeScript (Next.js)	Rich admin UX, SSR performance, component library to render MCP-UI (card/table/form).
9	Frontend — Mobile	React Native	SSE/WebSocket streaming, consent & policy UI, offline caches.

3.4 Reference Deployment Topology

Shows region distribution, cluster layout, and node pools for dev, staging, prod.

Diagram / table of regions, VNETs, node pools

3.5 Network & Security

Details ingress/egress, WAF, service-mesh, mTLS, and subnet isolation for zero-trust.

- Ingress paths, WAF, mTLS
- Zero-trust zones diagram

API and Data Models

Data Model :

1. skill_implementations

</> GraphQL

```
1 {
2   "_id": { "$oid": "66f101a0a1b2c30100000201" },
3   "skillId": "skill-astro-7c9e-4c3a-9c1a-1b2a3c4d5e6f",
4
5   "tenantId": "tenant-jio", // NEW for multi-tenant isolation
6   "userId": "ravindra.singh@ril.com",
7
8   "skillName": "Astrology",
9   "skillDescription": "Provides daily/weekly horoscopes, detailed birth charts
(Kundali), and relationship compatibility analysis.",
10
11   "partnerName": "AstroVision",
12   "partnerDescription": "AstroVision offers accurate astrology computations and curated
interpretations.",
13
14   "vectorSkillId": "vec-astro-7c9e-4c3a-9c1a-1b2a3c4d5e6f",
15
16   "prompts": { // NEW: agent construction at runtime
17     "system": "You are an expert Vedic astrology assistant. Be precise but empathetic.
Explain results in simple language and avoid deterministic predictions.",
18     "toolUseGuidelines": "If birth details are provided use get_birth_chart; otherwise
get_horoscope. For relationship queries call match_compatibility.",
19     "safety": "Never output exact birth time publicly; mask to nearest 5 minutes unless
user has consent."
20   },
21
22   "tags": ["astrology","horoscope","kundali","compatibility"], // NEW: improves SDS
retrieval & filters
23
24   "policy": { // NEW
25     "visibility": "public",
26     "rate_limits": { "rpm": 120 }
27   },
28
29   "mcpBindings": ["bind-astro-primary-v1"], // NEW: binds this skill to an MCP endpoint
via bindings
30
31   "specRefs": { // NEW: move big specs out of the doc
32     "openapi": {
33       "blobRef": "gs://partner-specs/astrovision/openapi@sha256:cafebabe1234",
34       "hash": "sha256:cafebabe1234",
35       "model": "3.0.1"
36     }
37   },
38
39   "createdBy": "ravindra.singh@ril.com",
40   "updatedBy": "ravindra.singh@ril.com",
41   "createdAt": { "$date": "2025-07-06T09:24:00Z" },
42   "updatedAt": { "$date": "2025-07-06T09:30:00Z" },
43
44   "status": "APPROVED",
45   "markForApproval": false,
46   "version": 3,
47   "approvalRemarks": "",
48   "_class": "com.rjil.adminPortal.model.SkillImplementation"
```

```

49  }
50
51
52

```

	Field	Type	Example / Allowed Values	Description & Significance
1	<code>_id</code>	ObjectId	66f101a0a1b2c30100000201	Mongo primary key.
2	<code>skillId</code>	string	skill-astro-7c9e-4c3a-9c1a-1b2a3c4d5e6f	Stable identifier used across services and in retrieval payloads.
3	<code>tenantId</code>	string	tenant-jio	Multi-tenant isolation; included in all reads/search filters.
4	<code>userId</code>	string (email/ID)	ravindra.singh@ril.com	Last editor/owner for audit.
5	<code>skillName</code>	string	Astrology	Display name; used in directory UI & text search.
6	<code>skillDescription</code>	string	Provides daily/weekly horoscopes...	Embedding source + directory description.
7	<code>partnerName</code>	string	AstroVision	Partner brand shown in portal.
8	<code>partnerDescription</code>	string	AstroVision offers accurate...	Longer partner bio for UI/context.
9	<code>vectorSkillId</code>	string	vec-astro-7c9e-...	Logical grouping for embeddings (useful for reindex & cleanup).
10	<code>prompts.system</code>	string	You are an expert Vedic astrology assistant...	Base system prompt to construct the agent.
11	<code>prompts.toolUseGuidelines</code>	string	If birth details are provided use get_birth_chart...	Sequencing/selection hints that improve tool reliability.
12	<code>prompts.safety</code>	string	Never output exact birth time publicly...	Guardrails, masking, and policy reminders.

13	<code>tags</code>	<code>array<string></code>	<code>["astrology", "horoscope", "kundali", "compatibility"]</code>	Improves Skill Directory filtering & ranking.
14	<code>policy.visibility</code>	<code>enum</code>	<code>public tenant private</code>	Controls who can discover/use the skill via search.
15	<code>policy.rate_limits.rpm</code>	<code>integer</code>	<code>120</code>	Soft per-skill throttle hint; can inform routing/backpressure.
16	<code>mcpBindings</code>	<code>array<string></code>	<code>["bind-astro-primary-v1"]</code>	References into <code>skill_mcp_bindings</code> (endpoint, per-tool policies, runtime budgets).
17	<code>specRefs.openapi.blobRef</code>	<code>string (URI)</code>	<code>gs://partner-specs/astrovision/openapi@sha256:cafebabe1234</code>	Pointer to large OpenAPI spec stored out-of-band.
18	<code>specRefs.openapi.hash</code>	<code>string</code>	<code>sha256:cafebabe1234</code>	Integrity check for the referenced spec blob.
19	<code>specRefs.openapi.model</code>	<code>string</code>	<code>3.0.1</code>	OpenAPI version of the referenced spec.
20	<code>createdBy</code>	<code>string (email/ID)</code>	<code>ravindra.singh@ril.com</code>	Audit field.
21	<code>updatedBy</code>	<code>string (email/ID)</code>	<code>ravindra.singh@ril.com</code>	Audit field.
22	<code>createdAt</code>	<code>ISO datetime</code>	<code>2025-07-06T09:24:00Z</code>	Audit field.
23	<code>updatedAt</code>	<code>ISO datetime</code>	<code>2025-07-06T09:30:00Z</code>	Audit field.
24	<code>status</code>	<code>enum</code>	<code>DRAFT PENDING_APPROVAL APPROVED REJECTED</code> (example shown: <code>APPROVED</code>)	Workflow state; used to gate publish and search exposure.
25	<code>markForApproval</code>	<code>boolean</code>	<code>false</code>	Signals item needs reviewer attention.

26	version	integer	3	Incremented on updates; helpful for review diffs.
27	approvalRemarks	string	""	Approver/reviewer comments.
28	_class	string (optional)	com.rjil.adminPorta l.model.SkillImpleme ntation	Framework/ORM class name (legacy/optional).

2. mcp_endpoints

</> GraphQL

```
1 {
2   "_id": { "$oid": "66f101a0a1b2c30100000101" },
3   "endpointId": "mcp-astrovision-primary",
4
5   "tenantId": "tenant-jio",
6   "partnerId": "astrovision",
7
8   "name": "AstroVision MCP",
9   "baseUrl": "wss://mcp.api.astrovision.com/mcp",
10  "transport": "ws",
11  "protocolVersion": "mcp-1.0",
12  "serverVersion": "1.0.0",
13
14  "auth": {
15    "type": "bearer",
16    "secret_ref": "secrets/partners/astrovision/mcp_token"    // never store raw secret
17  },
18
19  "capabilities": {
20    "tools": true,
21    "resources": false,
22    "prompts": false,
23    "ui": {                                // NEW (optional MCP UI support)
24      "supported": true,
25      "spec": "mcp-ui-0.1",
26      "components": ["card", "table", "form"],
27      "maxPayloadKb": 64
28    }
29  },
30
31  "discoveredTools": [                    // captured by Validation Service
32    { "name": "get_horoscope", "schemaHash": "sha256:5f2c..." },
33    { "name": "get_birth_chart", "schemaHash": "sha256:ab44..." },
34    { "name": "match_compatibility", "schemaHash": "sha256:d1e9..." }
35  ],
36  "discoveredUi": [
37    { "component": "card", "schemaHash": "sha256:1111...", "lastSeenAt": "2025-09-05T06:10:00Z" },
38    { "component": "form", "schemaHash": "sha256:2222..." }
39  ],
40
41  "health": {
42    "status": "up",
43    "avgLatencyMs": 210,
44    "lastChecked": "2025-07-06T09:35:00Z",
45    "uptime7dPct": 99.7
46  },
47
48  "createdAt": "2025-07-06T09:20:00Z",
49  "updatedAt": "2025-07-06T09:35:00Z"
50 }
51
52
53
```

	Field	Type	Example / Allowed Values	Description & Significance
1	<code>_id</code>	ObjectId	<code>66f101a0a1b2c30100000101</code>	Mongo primary key. [System-managed]
2	<code>endpointId</code>	string	<code>mcp-astrovision-primary</code>	Stable endpoint identifier referenced by bindings (<code>endpointRef</code>). [System-managed] (generated), visible in portal.
3	<code>tenantId</code>	string	<code>tenant-jio</code>	Multi-tenant isolation; requests must match tenant. [Partner input] (or prefilled by portal context).
4	<code>partnerId</code>	string	<code>astrovision</code>	Groups endpoints by partner for RBAC, reporting, quotas. [Partner input] (selected from partner org).
5	<code>name</code>	string	<code>AstroVision MCP</code>	Human-readable label for ops/portal. [Partner input]
6	<code>baseUrl</code>	URL	<code>wss://mcp.api.astrovision.com/mcp</code>	MCP connection address (ws/http). Critical for dialing. [Partner input]
7	<code>transport</code>	enum	<code>ws</code> <code>http</code> <code>stdio</code>	Selects MCP client/transport. [Partner input]
8	<code>protocolVersion</code>	string	<code>mcp-1.0</code>	MCP protocol claimed by server; we verify on handshake. [Partner input] , validated [Auto-discovered]
9	<code>serverVersion</code>	string	<code>1.0.0</code>	Vendor build/semver shown in ops. [Partner input] (optional), may be confirmed [Auto-discovered]
10	<code>auth.type</code>	enum	<code>bearer</code> <code>oauth2</code> <code>mtls</code>	Auth scheme to use. [Partner input]

11	<code>auth.secret_ref</code>	string	<code>secrets/partners/as trovision/mcp_token</code>	Reference to token/secret in Vault/KMS. We never store plaintext. [Partner input] (captured via secure form; stored as ref).
12	<code>auth.oauth2.token_url</code>	URL	<code>https://idp.example.com/oauth/token</code>	OAuth2 token endpoint (if <code>oauth2</code>). [Partner input]
13	<code>auth.oauth2.client_id_ref</code>	string	<code>secrets/.../client_id</code>	Secret ref for OAuth2 client id. [Partner input]
14	<code>auth.oauth2.client_secret_ref</code>	string	<code>secrets/.../client_secret</code>	Secret ref for OAuth2 client secret. [Partner input]
15	<code>auth.oauth2.scopes</code>	array<string>	<code>["mcp.read", "mcp.write"]</code>	Requested OAuth scopes. [Partner input]
16	<code>capabilities.tools</code>	boolean	<code>true</code>	Server implements MCP tools. [Partner input] (declared) and verified [Auto-discovered] via <code>list_tools</code> .
17	<code>capabilities.resources</code>	boolean	<code>false</code>	Exposes MCP resource catalog. [Partner input] (optional) and verified [Auto-discovered]
18	<code>capabilities.prompts</code>	boolean	<code>false</code>	Exposes MCP prompts catalog. [Partner input] (optional) and verified [Auto-discovered]
19	<code>capabilities.ui.supported</code>	boolean	<code>true</code>	Server can emit UI frames (MCP-UI extension). [Partner input] (optional) and probed [Auto-discovered]
20	<code>capabilities.ui.spec</code>	string	<code>mcp-ui-0.1</code>	UI spec version the client understands. [Partner input] (optional)
21	<code>capabilities.ui.components</code>	array<enum>	<code>["card", "table", "form"]</code>	Component types the server may emit. [Partner input] (optional), validated [Auto-discovered]
22	<code>capabilities.ui.maxPayloadKb</code>	integer	<code>64</code>	Payload cap to enforce client-side. [Partner input] (optional, policy hint)

23	<code>discoveredTools[].name</code>	string	<code>get_horoscope</code>	Tool name from <code>list_tools</code> . Source of truth for available tools. [Auto-discovered] (Validation writes it).
24	<code>discoveredTools[].schemaHash</code>	string	<code>sha256:5f2c...</code>	Hash of input schema for drift detection. [Auto-discovered]
25	<code>discoveredUi[].component</code>	enum	<code>card table form</code>	UI component observed during validation. [Auto-discovered]
26	<code>discoveredUi[].schemaHash</code>	string	<code>sha256:1111...</code>	UI schema hash for drift/change detection. [Auto-discovered]
27	<code>discoveredUi[].lastSeenAt</code>	ISO datetime	<code>2025-09-05T06:10:00Z</code>	Last time this UI component was seen. [System-managed] (set by Validation Service).
28	<code>health.status</code>	enum	<code>up degraded down</code>	Latest health from probe; used for routing/boosts. [Auto-discovered] and [System-managed]
29	<code>health.avgLatencyMs</code>	number	<code>210</code>	Moving average latency for endpoint. [Auto-discovered] and [System-managed]
30	<code>health.lastChecked</code>	ISO datetime	<code>2025-07-06T09:35:00Z</code>	Freshness of health signal. [System-managed]
31	<code>health.uptime7dPct</code>	number	<code>99.7</code>	7-day availability; can boost ranking. [System-managed]
32	<code>createdAt</code>	ISO datetime	<code>2025-07-06T09:20:00Z</code>	Audit field. [System-managed]
33	<code>updatedAt</code>	ISO datetime	<code>2025-07-06T09:35:00Z</code>	Audit field. [System-managed]

3. `skill_mcp_bindings` (how the skill uses the endpoint: policies + runtime + UI)

</> GraphQL

```

1 {
2   "_id": { "$oid": "66f101a0a1b2c30100000301" },
3   "bindingId": "bind-astro-primary-v1",
4
5   "tenantId": "tenant-jio",
6   "skillId": "skill-astro-7c9e-4c3a-9c1a-1b2a3c4d5e6f",
7   "skillVersion": "1.0.0",           // pin to the prompts/policies set
8 }

```

```
9  "endpointRef": "mcp-astrovision-primary",
10
11  "toolsAllowlist": ["get_horoscope","get_birth_chart","match_compatibility"],
12
13  "toolPolicies": {
14    "get_horoscope": {
15      "timeoutMs": 4000,
16      "rateLimit": { "rpm": 200 },
17      "paramDefaults": { "timezone": "Asia/Kolkata", "lang": "en" },
18      "paramRequirements": ["sign","period"],
19      "privacy": { "maskBirthTimeToNearestMinutes": 5 },
20      "source": { // optional provenance metadata
21        "timeoutMs": "server:x-timeout-ms",
22        "rateLimit": "server:x-rate-limit-rpm",
23        "paramDefaults": "schema.properties.*.default",
24        "paramRequirements": "schema.required"
25      },
26      "lastSyncedAt": "2025-07-06T09:33:00Z"
27    },
28    "get_birth_chart": {
29      "timeoutMs": 8000,
30      "rateLimit": { "rpm": 120 },
31      "paramRequirements": ["dob","tob","lat","lon"]
32    },
33    "match_compatibility": {
34      "timeoutMs": 7000,
35      "rateLimit": { "rpm": 100 },
36      "paramDefaults": { "model": "vedic" }
37    }
38  },
39
40  "runtime": { // agent runtime budgets for this binding
41    "maxSteps": 4,
42    "maxTokens": 2048,
43    "temperature": 0.1
44  },
45
46  "uiAllowlist": ["card","table","form"], // optional, if MCP UI is enabled
47  "uiPolicies": {
48    "maxInteractiveSteps": 3,
49    "allowExternalLinks": false,
50    "sanitizeHtml": true,
51    "maxPayloadKb": 64
52  },
53
54  "validation": {
55    "status": "passed",
56    "lastRunId": "val-astro-2025-07-06-001",
57    "latencyP95Ms": 620
58  },
59
60  "createdAt": "2025-07-06T09:26:00Z",
61  "updatedAt": "2025-07-06T09:34:00Z"
62 }
63
64
65
```

	Field	Type	Example / Allowed Values	Description & Significance
1	<code>_id</code>	ObjectId	<code>66f101a0a1b2c30100000301</code>	Mongo primary key. [System-managed]
2	<code>bindingId</code>	string	<code>bind-astro-primary-v1</code>	Stable join key referenced by skills and bootstrap API. [System-managed] (generated)

3	<code>tenantId</code>	string	<code>tenant-jio</code>	Tenant isolation for runtime reads and edits. [Partner input] (or set by portal context)
4	<code>skillId</code>	string	<code>skill-astro-7c9e-...</code>	Parent skill this binding belongs to. [System-managed] (derived when creating binding)
5	<code>skillVersion</code>	string	<code>1.0.0</code>	Pins prompt/policy snapshot used at runtime. [Partner input] (or selected during publish)
6	<code>endpointRef</code>	string	<code>mcp-astrovision-primary</code>	References <code>mcp_endpoints.endpointId</code> to attach the server. [Partner input] (chosen from endpoints)
7	<code>toolsAllowlist</code>	array<string>	<code>["get_horoscope", "get_birth_chart", "match_compatibility"]</code>	Only these tools can be called by the agent. [Partner input] (informed by discovered tools)
8	<code>toolPolicies.<tool>.timeoutMs</code>	integer (ms)	<code>4000</code>	Per-tool timeout budget; agent enforces. [Partner input] (may be prefilled from server hints)
9	<code>toolPolicies.<tool>.rateLimit.rpm</code>	integer	<code>200</code>	Per-tool RPM; enforced by agent/gateway. [Partner input] (may be prefilled)
10	<code>toolPolicies.<tool>.paramDefaults</code>	object	<code>{ "timezone": "Asia/Kolkata", "lang": "en" }</code>	Defaults injected when inputs omit them. [Partner input] (prefilled from schema defaults [Auto-discovered])
11	<code>toolPolicies.<tool>.paramRequirements</code>	array<string>	<code>["sign", "period"]</code>	Required params; validation before tool call. [Auto-discovered] from JSON Schema <code>required</code> , editable by partner
12	<code>toolPolicies.<tool>.privacy</code>	object	<code>{ "maskBirthTimeToNearEstMinutes": 5 }</code>	Privacy/PII rules applied in agent. [Partner input] (may be hinted via <code>x-privacy</code> [Auto-discovered])

13	<code>toolPolicies.</code> <code><tool>.source</code>	object	<code>{</code> <code>"timeoutMs":"server:</code> <code>x-timeout-ms",</code> <code>"rateLimit":"server:</code> <code>x-rate-limit-rpm",</code> <code>"paramDefaults":"sch</code> <code>ema.properties.*.def</code> <code>ault",</code> <code>"paramRequirements":</code> <code>"schema.required" }</code>	Provenance of policy values (server hints vs schema vs manual). [System-managed] (written by Validation/merge)
14	<code>toolPolicies.</code> <code><tool>.lastSyncedAt</code>	ISO datetime	<code>2025-07-</code> <code>06T09:33:00Z</code>	When policy was last auto-populated from server hints. [System-managed]
15	<code>runtime.maxSteps</code>	integer	<code>4</code>	Max plan→act loops for the agent. [Partner input] (platform-owned policy)
16	<code>runtime.maxTokens</code>	integer	<code>2048</code>	Generation budget. [Partner input] (platform-owned policy)
17	<code>runtime.temperature</code>	number	<code>0.1</code>	Determinism vs creativity. [Partner input] (platform-owned policy)
18	<code>uiAllowlist</code> <i>(optional)</i>	array<enum>	<code>["card", "table", "fo</code> <code>rm"]</code>	UI components allowed (if MCP-UI enabled). [Partner input]
19	<code>uiPolicies.maxInter</code> <code>activeSteps</code> <i>(optional)</i>	integer	<code>3</code>	Cap UI back-and-forth loops per turn. [Partner input]
20	<code>uiPolicies.allowExt</code> <code>ernalLinks</code> <i>(optional)</i>	boolean	<code>false</code>	Block off-domain links in UI frames. [Partner input]
21	<code>uiPolicies.sanitize</code> <code>Html</code> <i>(optional)</i>	boolean	<code>true</code>	Enforce HTML sanitization on client. [Partner input]
22	<code>uiPolicies.maxPaylo</code> <code>adKb</code> <i>(optional)</i>	integer	<code>64</code>	Max UI payload size per message. [Partner input]
23	<code>validation.status</code>	enum	<code>passed</code> <code>failed</code> <code>pending</code>	Gate for publish/runtime. [System-managed] (set by Validation Service)
24	<code>validation.lastRunI</code> <code>d</code>	string	<code>val-astro-2025-07-</code> <code>06-001</code>	Link to latest validation run. [System-managed]

25	<div>validation.latencyP95Ms</div>	number	<div>620</div>	Performance SLO marker for this binding. [System-managed]
26	<div>createdAt</div>	ISO datetime	<div>2025-07-06T09:26:00Z</div>	Audit. [System-managed]
27	<div>updatedAt</div>	ISO datetime	<div>2025-07-06T09:34:00Z</div>	Audit. [System-managed]

4. validation_runs (immutable audit of checks)

</> GraphQL

```
1 {
2   "id": { "$oid": "66f101a0a1b2c30100000401" },
3   "runId": "val-astro-2025-07-06-001",
4
5   "tenantId": "tenant-jio",
6   "skillId": "skill-astro-7c9e-4c3a-9cla-1b2a3c4d5e6f",
7   "skillVersion": "1.0.0",
8   "endpointRef": "mcp-astrovision-primary",
9
10  "status": "passed",
11
12  "checks": [
13    { "check": "handshake", "ok": true, "latencyMs": 120 },
14    { "check": "listTools", "ok": true, "found": 3 },
15    { "check": "toolSchema:get_horoscope", "ok": true },
16    { "check": "toolSmoke:get_horoscope", "ok": true, "latencyMs": 260 },
17    { "check": "ui:card", "ok": true, "latencyMs": 48 },           // optional, if
    UI probed
18    { "check": "authScope", "ok": true }
19  ],
20
21  "metrics": { "latencyP95Ms": 620, "latencyAvgMs": 310 },
22  "errors": [],
23  "langfuseTraceId": "lf-astro-001",
24
25  "createdAt": "2025-07-06T09:32:00Z"
26 }
27
28
29
```

	Field	Type	Example / Allowed Values	Description & Significance
1	<div>_id</div>	ObjectId	<div>66f101a0a1b2c30100000401</div>	Mongo primary key. [System-managed]
2	<div>runId</div>	string	<div>val-astro-2025-07-06-001</div>	Stable validation job id used for audit and joins from bindings. [System-managed] (generated by Validation Service)

3	<code>tenantId</code>	string	<code>tenant-jio</code>	Tenant isolation for validation context. [System-managed] (from portal/session)
4	<code>skillId</code>	string	<code>skill-astro-7c9e-4c3a-9c1a-1b2a3c4d5e6f</code>	Skill under validation. [System-managed] (resolved from request)
5	<code>skillVersion</code>	string	<code>1.0.0</code>	Version being validated (draft or published). [System-managed] (derived when job starts; partner may choose version in UI, but record is written by system)
6	<code>endpointRef</code>	string	<code>mcp-astrovision-primary</code>	MCP endpoint tested. [System-managed] (copied from binding at job start)
7	<code>status</code>	enum	<code>passed</code> <code>failed</code> <code>pending</code>	Overall outcome of the validation job. [System-managed]
8	<code>checks[].check</code>	string	<code>handshake</code> , <code>listTools</code> , <code>toolSchema:get_horo</code> <code>scope</code> , <code>toolSmoke:get_horos</code> <code>cope</code> , <code>ui:card</code> , <code>authScope</code>	Name of each probe executed (handshake, discovery, schema verify, smoke calls, optional UI probe). [System-managed] with results [Auto-discovered]
9	<code>checks[].ok</code>	boolean	<code>true/false</code>	Pass/fail per probe. [System-managed]
10	<code>checks[].latencyMs</code> (optional)	number	<code>120</code>	Latency of a probe (ms). [System-managed]
11	<code>checks[].found</code> (optional)	number	<code>3</code>	Count metric for a probe (e.g., tools discovered). [System-managed]
12	<code>checks[].details</code> (optional)	object	<code>{ "tools": ["get_horoscope", "..."] }</code>	Small structured payload for extra context (kept lightweight). [System-managed] with values [Auto-discovered]
13	<code>metrics.latencyP95Ms</code>	number	<code>620</code>	Aggregated p95 latency across smoke calls (SLO signal). [System-managed]

14	<code>metrics.latencyAvgMs</code>	number	<code>310</code>	Aggregated average latency across probes. [System-managed]
15	<code>errors[]</code>	array<string>	<code>["auth failed: 401 during handshake"]</code>	Top-level errors, if any, for triage. [System-managed]
16	<code>langfuseTraceId</code> <i>(optional)</i>	string	<code>lf-astro-001</code>	Correlates the run with traces/spans in Langfuse. [System-managed]
17	<code>createdAt</code>	ISO datetime	<code>2025-07-06T09:32:00Z</code>	When the job record was created. [System-managed]

5. `sds_index_map` (lets SDS delete/rollback by exact IDs)

</> GraphQL

```
1 {
2   "_id": { "$oid": "66f101a0a1b2c30100000601" },
3
4   "tenantId": "tenant-jio",
5   "skillId": "skill-astro-7c9e-4c3a-9c1a-1b2a3c4d5e6f",
6   "versionId": "sv-astro-1-0-0",
7   "environment": "prod",           // "prod" | "sandbox" (if/when sandbox returns)
8
9   "collection": "skills_v1",
10  "pointIds": [1200101, 1200102, 1200103],
11  "numPoints": 3,
12
13  "model": "text-embedding-3-large",
14  "chunking": { "algo": "recursive", "maxTokens": 512, "overlap": 64 },
15
16  "checksum": { "textHash": "sha256:...", "configHash": "sha256:..." },
17  "batchId": "idx-2025-09-04-astro-0001",
18  "active": true,
19
20  "createdAt": "2025-09-04T10:30:00Z",
21  "updatedAt": "2025-09-04T10:31:00Z"
22 }
23
24
25
```

	Field	Type	Example / Allowed Values	Description & Significance
1	<code>_id</code>	ObjectId	<code>66f101a0a1b2c30100000601</code>	Mongo primary key. [System-managed]
2	<code>tenantId</code>	string	<code>tenant-jio</code>	Tenant isolation for vector ownership and deletes. [System-managed] (copied from skill)

3	<code>skillId</code>	string	<code>skill-astro-7c9e-4c3a-9c1a-1b2a3c4d5e6f</code>	Skill whose content was embedded. [System-managed]
4	<code>versionId</code>	string	<code>sv-astro-1-0-0</code>	Exact skill version for these vectors (enables rollback). [System-managed]
5	<code>environment</code>	enum	<code>prod</code> <code>sandbox</code>	Indexing environment/namespace. <code>sandbox</code> used for draft testing. [System-managed]
6	<code>collection</code>	string	<code>skills_v1</code> (or <code>skills_v1_sandbox</code>)	Qdrant collection name where vectors were written. [System-managed]
7	<code>pointIds</code>	array<int or string>	<code>[1200101, 1200102, 1200103]</code>	Exact Qdrant point IDs for this version. Enables deterministic delete/rollback. [Auto-discovered] (returned by Qdrant on upsert)
8	<code>numPoints</code>	integer	<code>3</code>	Count of vectors written for this version. [System-managed]
9	<code>model</code>	string	<code>text-embedding-3-large</code>	Embedding model used (audit, reproducibility). [System-managed] (from SDS config)
10	<code>chunking.algo</code>	string	<code>recursive</code>	Chunking strategy used. [System-managed] (from SDS config)
11	<code>chunking.maxTokens</code>	integer	<code>512</code>	Max tokens per chunk. [System-managed]
12	<code>chunking.overlap</code>	integer	<code>64</code>	Token overlap between chunks. [System-managed]
13	<code>checksum.textHash</code>	string	<code>sha256:...</code>	Hash of concatenated indexed text; detects drift/regression. [System-managed]
14	<code>checksum.configHash</code>	string	<code>sha256:...</code>	Hash of model/chunking config; ensures reproducibility. [System-managed]
15	<code>batchId</code>	string	<code>idx-2025-09-04-astro-0001</code>	Indexing batch/job identifier (join to logs). [System-managed]

16	<code>active</code>	boolean	<code>true/false</code>	<code>true</code> while version is live; set <code>false</code> on rollback/delete. [System-managed]
17	<code>createdAt</code>	ISO datetime	<code>2025-09-04T10:30:00Z</code>	When indexing map was created. [System-managed]
18	<code>updatedAt</code>	ISO datetime	<code>2025-09-04T10:31:00Z</code>	Last update time (e.g., reindex). [System-managed]
19	<code>expiresAt</code> (optional)	ISO datetime	<code>2025-09-19T10:31:00Z</code>	TTL for sandbox entries; rarely used in prod. [System-managed]

6. `skill_versions` (optional, if in case we keep version docs)

`</> GraphQL`

```
1
2 {
3   "_id": { "$oid": "66f101a0a1b2c30100000501" },
4   "versionId": "sv-astro-1-0-0",
5   "skillId": "skill-astro-7c9e-4c3a-9c1a-1b2a3c4d5e6f",
6
7   "semver": "1.0.0",
8   "changelog": "Initial public release with horoscope, birth chart, and compatibility.",
9   "publishedBy": "ravindra.singh@ril.com",
10
11  "status": "published",
12  "embedding": { "model": "text-embedding-3-large", "point_ids": [1200101, 1200102] },
13  "metadata": { "retrieval_weight": 1.0 },
14
15  "createdAt": "2025-07-06T09:28:00Z"
16 }
17
18
```

Indexes :

`</> GraphQL`

```
1 // skill_implementations
2 db.skill_implementations.createIndex({ tenantId: 1, skillId: 1 }, { unique: true });
3 db.skill_implementations.createIndex({ tenantId: 1, skillName: 1 });
4 db.skill_implementations.createIndex({ skillName: "text", skillDescription: "text",
5   tags: "text" });
6
7 // mcp_endpoints
8 db.mcp_endpoints.createIndex({ endpointId: 1 }, { unique: true });
9 db.mcp_endpoints.createIndex({ tenantId: 1, partnerId: 1 });
10 db.mcp_endpoints.createIndex({ "health.status": 1, updatedAt: -1 });
11
12 // skill_mcp_bindings
13 db.skill_mcp_bindings.createIndex({ skillId: 1, skillVersion: 1 });
14 db.skill_mcp_bindings.createIndex({ endpointRef: 1 });
15 db.skill_mcp_bindings.createIndex({ "validation.status": 1 });
16
17 // validation_runs
18 db.validation_runs.createIndex({ runId: 1 }, { unique: true });
19 db.validation_runs.createIndex({ skillId: 1, skillVersion: 1, endpointRef: 1, createdAt:
20   -1 });
```

```
19
20 // sds_index_map
21 db.sds_index_map.createIndex({ tenantId: 1, skillId: 1, versionId: 1 }, { unique: true
  });
22 db.sds_index_map.createIndex({ environment: 1, active: 1 });
23
24
25
```

This section covers all internal/external APIs. Define interfaces and Data models (DB Schemas/internal POJO/ DTO / DAO / JSON etc...)

System Flow Summary

Step-by-step narrative of a request—from user input through orchestration to response. Captured above in Key Sequence section

Observability & Instrumentation

States what metrics, logs, and traces are collected and how alerts are triggered.

	≡ Signal	≡ Tool	≡ Notes
1	Logs	Databricks Dashboard	Structured logs with correlation IDs: ts, level, service, requestId, tenantId, skillId, bindingId, endpointRef, event (e.g., MCP_CALL_START, SDS_SEARCH, BOOTSTRAP_READ). Keep payloads summarized; mask PII; secrets never logged. Retention: 7d hot, 30–90d warm (configurable).
2	Alerts	Alertmanager → PagerDuty	SLO/SLA-driven paging with runbooks. Examples: (1) SDS p95 > 100ms for 10m; (2) Orchestrator 5xx > 1% for 5m; (3) MCP endpoint health = down 5m; (4) Validation failure rate > 10% over 30m; (5) Kafka lag > 10k for 10m; (6) Qdrant error rate > 0.5% for 5m; (7) Vault secret fetch failures > 5/min; (8) LLM cost burn > daily budget. Alerts include Grafana + Langfuse deep links.

Assumptions & Constraints

Lists what the design presumes to be true (e.g., cloud region, budget) and any hard limits (e.g., data-residency).

Assumptions

- **Identity & Secrets**
 - OIDC (tenant-aware) for users; JWT/mTLS for services.
 - Partners provide credentials once; the portal stores **only** `secret_ref` (no plaintext secrets in DB).
- **Protocol / Interop**
 - MCP protocol `mcp-1.0` is used by partners; servers implement `list_tools` and tool invocation consistently.
 - Optional **MCP-UI** frames may be emitted by some partners; clients support a safe subset (card/table/form).
- **Validation Gate**
 - Skills/bindings must **pass validation** before they are eligible for publish and runtime (`strict=true`).
- **Time & Locale**
 - Timestamps stored in UTC; examples use **Asia/Kolkata**; localization handled at the client/UI layer.
- **Sandbox**
 - Sandbox/testing workflow is **deferred**; current scope focuses on production runtime.

Constraints

- **Security & Compliance**
 - **No raw secrets** stored anywhere—only `secret_ref` /KMS references. Token resolution happens at runtime.
 - PII/inputs to tools are **redacted/summarized** in traces; no chain-of-thought persisted.
- **Access Control**
 - All reads/writes are **tenant-scoped**; cross-tenant data access is prohibited.
 - Only **allowlisted tools** in a binding can be invoked at runtime.
- **Validation & Publication**
 - Binding `validation.status` must be `passed` for bootstrap/use; otherwise service returns **409** (when `strict=true`).
 - Tool policy provenance is recorded; platform values **override** partner hints.
- **Capacity & Limits (guardrails)**
 - **MCP-UI frames**: only **allowlisted components**, payload ≤ 64 KB, **no external links** unless allowed by policy.
- **Change Management**
 - Publishing creates a new `versionId`; SDS reindex is **idempotent**; rollback uses `sds_index_map`.
 - Backward-incompatible schema changes require versioned APIs or migrations before rollout.
 - **Third-Party Dependencies**
 - MCP endpoints must provide stable SLAs and notify on schema changes; breaking changes without notice may lead to **temporary unlisting**.

Security Considerations

Catalogues auth, RBAC, malware scanning, audit logging, and key threat mitigations across the Partner Portal → Validation → SDS → Runtime (Orchestrator & Skill Agent Service) path.

1) Identity, Auth & RBAC (Need to discuss further)

- **User AuthN (Portal)**: OAuth2 / OIDC via JioID issuing JWT (short-lived access tokens, refresh via PKCE).
- **Service AuthN**: mTLS between services; JWT for intra-API calls with `aud/iss/exp` checks and key rotation.
- **AuthZ (RBAC)**: Tenant-scoped roles — **Partner Owner / Editor / Viewer** — plus Platform roles (**Ops, Validator**).
- **Skill invocation RBAC**: Only published skills with `policy.visibility` permitted; `binding.validation.status=passed` required when `strict=true`.
- **ABAC**: Enforce `tenantId` on every read/write; **cross-tenant access is denied** by default.
- **Principle of Least Privilege**: SAS only receives the **binding** + `secret_ref` it needs for the chosen skill.

2) Secrets & Keys

- **Secret storage:** No plaintext secrets in Mongo; store `secret_ref` only. Backed by Google Secrets Manager.
- **Retrieval:** SAS resolves secret refs **at call time**; results held in memory only; **never logged**.
- **Rotation:** Partner tokens rotated per policy (e.g., 90 days). Validation re-runs on rotation.
- **Key management:** KMS-managed keys for envelope encryption (optional field-level encryption for sensitive prompts).

3) Data Protection (At Rest & In Transit)

- **Transport:** TLS 1.2+ externally; mTLS internally. WS (MCP), HTTP, and gRPC all terminate at gateway/ingress with WAF.
- **At rest:** Disk encryption for Mongo/Qdrant; object storage for specs is KMS-encrypted.
- **Field-level protections:** Optional encryption for prompt fields; PII masking per policy.

4) Malware & Content Scanning

- **Uploads:** Any partner-uploaded files (e.g., OpenAPI specs, sample payloads) are **scanned before LLM access** and before persisting references.
- **Runtime attachments:** If MCP or tools return file links, enforce **allowlisted MIME types**, size caps, and scan before exposing to clients.
- **Blocking:** Quarantine & reject on detection; notify partner and log an audit event.

5) Input/Output Safety (LLM + Tools)

- **Prompt-injection hardening:**
 - Route only to **allowlisted tools** (`toolsAllowlist`), **never** open-ended code exec.
 - **Param validation** against tool JSON Schemas; reject unknown/extra keys.
 - Budget guardrails from binding `runtime` (maxSteps, tokens, timeoutMs per tool).
- **Data exfil prevention:**
 - No secrets in prompts/responses; redact credentials/PII in traces and logs.
 - SAS **egress restriction** to declared MCP `baseURL` (DNS pin / IP allowlist) to mitigate SSRF.
- **MCP-UI safety (if enabled):**
 - Enforce `uiAllowlist` and `uiPolicies` (payload ≤ N KB, **no external links** unless explicitly allowed, strict HTML sanitization).
 - Clients ignore unknown components; degrade to text.

6) Network, Platform & Supply Chain

- **Gateway/WAF:** Rate limiting, IP allowlists, header normalization, request/response size limits, bot/DDoS protections.
- **Kubernetes hardening:** Non-root containers, read-only root FS, seccomp/AppArmor, NetworkPolicies (east-west), image signing (Cosign), SBOM + SCA scans.
- **Dependency hygiene:** Automated vulnerability scans; critical CVEs block deploy until patched.
- **Change control:** GitOps (Argo CD) with reviews; infrastructure as code (Terraform).
- **Multi-region:** Data residency honored per tenant; geo-fencing of compute/DB if required.

7) Logging, Tracing & Audit

- **Structured logs:** PII-safe JSON with correlation IDs (`requestId`, `tenantId`, `skillId`, `bindingId`, `endpointRef`).
- **Traces:** OpenTelemetry + **Langfuse** with redacted attributes (no raw tool inputs, no secrets).
- **Audit trail:** Append-only events for **who changed prompts/bindings**, **who validated/published**, **who invoked** which skill/tool.
- **Tamper resistance:** Write-once storage tier (or hash-chaining) for sensitive audit records; time sync via NTP.

8) Alerts & Incident Response (IR)

- **Alerting:** PagerDuty on SLO/RBA violations (auth anomalies, MCP health down, SDS latency spikes, secret fetch failures).

- **IR playbooks:** Threat-specific runbooks (token leak, SSRF, prompt-injection incident, schema-drift rollback), with auto-capture of relevant traces/logs.
- **Access reviews:** Quarterly RBAC audits; break-glass accounts gated and logged.

9) Threats & Mitigations (Quick Map)

	☰ Threat	☰ Mitigations
1	Prompt injection / tool hijacking	Tool allowlist at binding; strict JSON Schema validation; runtime maxSteps/tokens ; strip/system-prompt hardening; no chain-of-thought logging.
2	Cross-tenant data leak	Tenant-scoped queries everywhere; SDS/Qdrant payload filters (<code>tenant_id</code>); authZ checks at gateway and service; per-tenant encryption keys (optional).
3	Credential exfiltration	Only <code>secret_ref</code> stored; secrets fetched just-in-time; redaction in logs/traces; egress pinned to MCP host; no reflective prompts.
4	SSRF / egress abuse	SAS egress restricted to allowlisted MCP baseUrl ; DNS pinning / IP allowlist; proxy with ACLs; block file://, 169.254/metadata.
5	Schema drift / breaking MCP changes	Validation job detects schemaHash changes; bindings require re-approval; route down-weights degraded endpoints automatically.
6	DoS / cost overrun	Rate limits (gateway + tool RPM); concurrency caps; LLM budgets; circuit breakers; back-pressure; alert on LLM cost burn .
7	Malware via uploads or MCP content	Pre-ingest AV scanning , MIME/type allowlists, size caps, quarantine + alert; disallow executable artifacts.
8	Malicious MCP-UI frames	<code>uiAllowlist</code> , strict sanitization, payload caps, disallow external links unless policy allows; client ignores unknown components.
9	MITM / token theft	TLS 1.2+ everywhere; mTLS internally; token bound to <code>aud/iss/exp</code> ; short-lived creds; Vault access via workload identity.
10	Unauthorized publish or policy change	RBAC + dual control for publish; audit log with immutable storage; notification on policy changes; require validation.status=passed .

10) Defaults & Guardrails (recommended)

- **JWT TTLs:** 15–30 min access, 24 h refresh; clock skew ≤ 60 s.

- **Rate limits:** Per-tenant global RPM + per-tool RPM from `toolPolicies`.
- **Payload limits:** Request/response \leq 1–5 MB (service-specific); MCP-UI \leq 64 KB (configurable).
- **Timeouts:** Tool `timeoutMs` defaults (e.g., 4–8s); Orchestrator E2E soft \leq 8s, hard \leq 20s.
- **Retention:** Logs 7d hot/30–90d warm; traces 30–90d; audit \geq 1 year (per policy).
- **Publish gate:** `validation.status=passed` + health not `down` in last 10 min.

Degradation & Fallback Strategy

Defines graceful-degradation tactics—**read-only mode**, **cached responses**, and **feature flags**

Modes & What Changes

	≡ Mode	≡ When it Activates	≡ What Changes (Scope)
1	Normal	All SLOs healthy	Full capabilities: SDS top-K retrieval, LLM routing, SAS planning, MCP tools (e.g., <code>get_horoscope</code> , <code>get_birth_chart</code>), optional MCP-UI.
2	Degraded (Routing/Tools)	MCP endpoint <code>health=down/degraded</code> or tool error-rate \uparrow ; SDS p95 > 120 ms for 10 min	Orchestrator re-routes to next healthy candidate; SAS limits planning (\leq <code>maxSteps_degraded</code>), reduces K, disables non-essential tools, suppresses MCP-UI frames.
3	Reduced-Cost Mode	Daily LLM budget \geq threshold or traffic surge	Switch to cheaper model, lower <code>maxTokens</code> , disable reranker/reasoning; enable aggressive caches.
4	Read-Only Mode (ROM)	Mongo/Vault writes failing, or publish queue backlog high	Portal: edits/publish disabled; SDS indexing paused (503); runtime reads allowed; services served from cache.
5	Maintenance / Safe Halt	Multiple critical deps failing	Serve cached apology/"status" responses; health endpoints only.

Trigger → Action Matrix

	≡ Trigger (from Observability)	≡ Automatic Actions	≡ Operator Playbook
1	MCP endpoint <code>down</code> 5 min or p95 > 800 ms	Orchestrator retries/backoff → pick next candidate; SAS disables affected tools; endpoint weight \downarrow	Notify partner; run validation; if persistent, temporarily unlist endpoint.

2	SDS p95 > 120 ms 10 min	Lower <code>retrieval.k</code> , enable candidate cache, pre-warm Redis; degrade reranker	Inspect Qdrant health/HNSW params, recent publish volume; scale SDS/Qdrant.
3	LLM cost burn > daily budget (Can be handled in LLM Router)	Switch <code>llm.model</code> to cost tier; cut <code>maxTokens</code> ; disable reasoning; cap concurrency	Review traffic; apply tenant-level quotas; tighten binding budgets.
4	Mongo write errors > 5%	Enter ROM	Investigate primary/replica health; failover; clear ROM after recovery.
5	Vault secret fetch failures > 5/min	Deny new MCP connects; serve from cached service	Check Vault health/permissions; rotate tokens if needed.
6	Validation failure rate > 10% 30 min	Block publish; keep runtime unaffected	Check schema drift; contact partners; roll back last binding changes.

Panic / Kill-Switch Mode

A hard stop designed to **contain risk instantly** (security incident, runaway cost, partner compromise) by halting risky paths, preserving forensics, and providing a safe user experience.

Objectives

- **Protect users & data** (immediately cut risky execution paths).
- **Bound blast radius** (scope by tenant/skill/endpoint/region or global).
- **Preserve evidence** (freeze audit/traces/logs for IR).
- **Offer clear UX** (status banners / static responses).
- **Enable controlled recovery** (gated re-enable with smoke checks).

Activation Scopes & Levels

	☰ Scope	☰ Examples	☰ Level	☰ What it means
1	Global	Platform-wide event (e.g., provider breach)	Halt	Stop all runtime invocations, MCP egress, publish/validate/index.
2	Region	Cloud zone outage	Halt / Read-only	Freeze writes in region; route reads to healthy region (if allowed).
3	Tenant	Single-tenant incident or legal hold	Halt / Read-only / Tool-less	Block that tenant's runtime & publish; optionally allow read-only portal.
4	Skill / Binding	Compromised prompts/policies, schema drift	Halt / Tool-less	Disable that skill or binding; other skills work.
5	Endpoint	Partner MCP compromised or unstable	Halt / Egress-block	Block SAS egress to that endpoint; re-route to alternates.

Compliance & Regulatory

Maps design choices to GDPR, RBI, ISO-27001, SOC-2, or other mandates.

- Data residency requirements
- Consent flows / audit log retention

Scalability & Extensibility

Describes autoscaling, sharding, and plug-in mechanisms for models, skills, tools.

	≡ Layer	≡ Strategy
1	Core services	Horizontal autoscale
2	Vector DB	Sharding per tenant
3	LLM Router	Pluggable providers

Future Roadmap

- Multi-modal extensions (video, AR)
- Edge LLM serving

Glossary

	≡ Acronym	≡ Meaning
1	MCP	Model Context Protocol
2	RAG	Retrieval-Augmented Generation
3	MyHub	JioOmni in-app user action center

Open Issues / Risks