

# Architecture Gap Analysis: SirsiNexusApp vs Strategic AI Roadmaps (SAR)

## 1. Sirsi Architecture Map

Sirsi is an "Enterprise-Grade AI Infrastructure Assistant" built with a polyglot, distributed-first mindset.

### Services & Modules

- **Core Engine (Rust)**: The "Hypervisor". Orchestrates internal agents, handles gRPC communication, and provides the main API gateway.
- **Connectors (Go)**: Low-level cloud provider integrations (AWS, GCP, Azure).
- **Planner (Python)**: High-level AI orchestration and planning logic.
- **Analytics Platform (Python)**: Forecasting, anomaly detection, and risk scoring using PyTorch/TensorFlow.
- **UI (Next.js/React)**: Desktop and mobile-ready web application.
- **CLI (Rust/Tauri)**: Local desktop client for infrastructure management.

### Datastore Inventory

- **CockroachDB**: Distributed SQL for resilient metadata and relational state.
- **Qdrant**: Vector database for infrastructure patterns and RAG support.
- **Redis**: High-speed cache and real-time context management.
- **Hedera**: Knowledge Graph distribution via DLT.

### Runtime Boundaries

- **Internal**: gRPC (port 50051) for high-performance agent-to-agent communication.
- **External**: WebSocket (port 8080) for real-time UI updates + REST (port 8081) for standard API.
- **Context**: Shared context managed via Redis and passed across agent boundaries.

### Deployment Topology

- **Local**: `docker-compose` for full-stack orchestration.
  - **Production**: Robust `k8s` (Kubernetes) manifests with Helm support, including auto-scaling, ingress (Nginx), and persistent storage.
-

## 2. SAR Architecture Map (Current)

SAR is a "Lean & Mean Multi-Tenant SaaS" optimized for rapid iteration on Vercel.

### Services & Modules

- **Backend (Express/TypeScript)**: Monolithic API handling authentication, multi-tenancy logic, and OpenAI Assistant integration.
- **Frontend (Vite/React)**: Lightweight SPA with complex state management for roadmap generation.
- **Shared Package**: Shared Zod schemas and utility types.

### Datastore Inventory

- **PostgreSQL (Neon)**: Serverless SQL via Neon, using Drizzle ORM for schema management.
- **AWS S3**: Artifact and document storage with presigned URLs.
- **OpenAI Assistant API**: Outsourced state/memory for agent conversations.

### Runtime Boundaries

- **Primary**: REST API over HTTPS.
- **LLM**: Server-side communication with OpenAI (stateful threads).

### Deployment Topology

- **Vercel**: Unified deployment for both Frontend and Backend (Serverless functions).
- 

## 3. Gap Table

Feature	SirsiNexusApp	Strategic AI Roadmaps (SAR)	Gap Assessment
<b>Backend Stack</b>	Polyglot (Rust, Python, Go)	Unified (TypeScript/Express)	Sirsi is significantly more complex but scalable; SAR is faster to modify.
<b>Service Boundaries</b>	Decentralized (gRPC microservices)	Monolithic (Express)	Sirsi handles multi-modal tasks (infra, analytics) as separate services.

Feature	SirsiNexusApp	Strategic AI Roadmaps (SAR)	Gap Assessment
<b>Persistence</b>	Distributed SQL (Cockroach) + Vector	Serverless SQL (Neon)	Sirsi is built for horizontal scale and high availability.
<b>Observability</b>	Prometheus, OpenTelemetry, Grafana	Minimal (Logging/Vercel Vitals)	Sirsi has enterprise-grade perf monitoring and anomaly detection.
<b>Infra/Deploy</b>	K8s + Helm + Docker Compose	Vercel (PaaS)	SAR relies on PaaS constraints; Sirsi owns the full stack.
<b>Real-time</b>	WebSocket + gRPC	Polling / Request-Response	Sirsi supports high-frequency interactive sessions.
<b>Data Protocol</b>	MCP (Model Context Protocol)	Custom Zod + JSON API	Sirsi uses standardized agent communication protocols.

## 4. "Borrow List" (Patterns for SAR)

### 1. Model Context Protocol (MCP):

- *Implementation:* Adopt MCP for SAR's background diagnostic generation to allow modular "Expert Agents" (e.g., SecurityExpert, BizModelExpert) to communicate via a standard interface.

### 2. Vector Store Integration (Qdrant Pattern):

- *Implementation:* Move from OpenAI's hidden vector search to a dedicated Qdrant or Pinecone instance to allow cross-tenant knowledge synthesis and better RAG control.

### 3. The "Hypervisor" Orchestrator:

- *Implementation:* Refactor the `backend/src/services` logic into an orchestrator service that manages agent transitions, rather than putting that logic in controllers.

### 4. Prometheus/OpenTelemetry Integration:

- *Implementation:* Integrate `prom-client` in the Express backend to track "Success Rate of AI Generations" and "Token Utilization" per tenant.

## 5. Schema-First Identity (Prisma-like):

- *Implementation:* While using Drizzle, adopt the Sirsi pattern of centralized schema packages (`packages/prisma`) to ensure the DB is the source of truth across multiple potential services.

## 6. Multi-Modal Testing (K6/Load-testing):

- *Implementation:* Adopt the `load-testing` folder pattern. Run K6 tests against the `/api/generate` endpoint to measure latency under load.
- 

## 5. Integration Option: "Safe Hardening Interface"

**Concept:** A minimal "Verification Gateway" where Sirsi provides security and reliability hardening for SAR without merging the two codebases.

### Proposal: The Sirsi-Nexus Hardening Proxy

- **Minimal Interface:** SAR exposes a specific "Hardening Webhook" or gRPC endpoint.
- **How it works:**
  1. When SAR generates a complex output (e.g., a "Strategic Implementation Plan"), it forwards the draft to a Sirsi-Nexus sidecar service.
  2. Sirsi-Nexus runs its **Compliance & Security Engine** (the `core-engine/src/compliance` logic) against the draft.
  3. Sirsi returns a "Hardening Report" (Score, Risks, Mitigations).
  4. SAR displays this report to the user as an "Enterprise Safety Badge."
- **Result:** SAR gains enterprise-grade validation logic without the operational overhead of running the full Sirsi stack.