

# RADIANT Platform - Administrator Guide

Complete guide for managing the RADIANT AI Platform via the Admin Dashboard

Version: 5.52.29 | Last Updated: January 2026

**Compliance Frameworks:** HIPAA, SOC 2 Type II, GDPR, FDA 21 CFR Part 11

---

## Table of Contents

1. Introduction
2. Accessing the Admin Dashboard
3. Dashboard Overview
4. Tenant Management
5. User & Administrator Management
6. AI Model Configuration
7. Provider Management
8. Billing & Subscriptions
9. Storage Management
10. Orchestration & Preference Engine
11. Pre-Prompt Learning
12. Localization
13. Configuration Management
14. Security & Compliance
15. Cost Analytics
16. Revenue Analytics
17. SaaS Metrics Dashboard
18. A/B Testing & Experiments
19. Audit & Monitoring
20. Database Migrations
21. API Management
22. Troubleshooting
23. Delight System Administration
24. Domain Ethics Registry
25. Model Proficiency Registry
26. Model Coordination Service
27. Ethics Pipeline
28. Inference Components
29. Service Environment Variables 31A. **Cato/Genesis Consciousness Architecture - Executive Summary**
30. Cato Global Consciousness Service
31. Cato Genesis System
32. AGI Brain - Project AWARE
33. Truth Engine™ - Project TRUTH
34. Advanced Cognition Services (v6.1.0)
35. Learning Architecture - Complete Overview 41B. Empiricism Loop (Consciousness Spark)
36. Genesis Cato Safety Architecture
37. Radiant CMS Think Tank Extension
38. AWS Free Tier Monitoring
39. Just Think Tank: Multi-Agent Architecture
40. RADIANT vs Frontier Models: Comparative Analysis
41. Flyte-Native State Management
42. Mission Control: Human-in-the-Loop (HITL) System
43. The Grimoire - Procedural Memory (NEW in v5.0)

44. The Economic Governor - Cost Optimization (NEW in v5.0)
  45. Self-Optimizing System Architecture (NEW in v5.0)
  46. Genesis Infrastructure: Sovereign Power Architecture
  47. Cato Security Grid: Native Network Defense
  48. AGI Brain & Identity Data Fabric: Agentic Orchestration
  49. Deployment Safety & Environment Management
  50. White-Label Invisibility (Moat #25)
  51. User Violation Enforcement
  52. Multi-Protocol Gateway Architecture
  53. RAWS v1.1 - Model Selection System
  54. Cortex Memory System
  55. Cortex Graph-RAG Knowledge Engine
  56. Complete Admin API Architecture
- 

## 1. Introduction

### 1.1 What is RADIANT?

RADIANT is a multi-tenant AWS SaaS platform providing unified access to 106+ AI models through:

- **50 External Provider Models:** OpenAI, Anthropic, Google, xAI, DeepSeek, and more
- **56 Self-Hosted Models:** Running on AWS SageMaker for cost control and privacy
- **Intelligent Routing:** Cognitive Router for optimal model selection
- **Preference Engine:** Personalization learning from user interactions

### 1.2 Administrator Roles

Role	Permissions	Use Case
<b>Super Admin</b>	Full access to all features	Platform owner
<b>Admin</b>	Tenant management, billing, models	Operations team
<b>Operator</b>	Read access, limited actions	Support team
<b>Auditor</b>	Read-only access to logs	Compliance team

**Role Details** **Super Admin** - The highest privilege level with unrestricted access:

- Create and delete tenants
- Manage all administrators
- Access all billing and financial data
- Modify system-wide configuration
- Approve production database migrations
- Impersonate any tenant for debugging
- Access compliance and audit reports
- Typically limited to 1-3 people (CTO, lead engineer)

**Admin** - Day-to-day operations management:

- Create and modify tenants (cannot delete)
- Manage users within tenants
- Configure AI models and providers
- View billing data (cannot modify pricing)
- Monitor system health
- Cannot access other admin accounts
- Typically assigned to operations team members

**Operator** - Limited support and monitoring:

- View tenant information (read-only)
- View user issues and support tickets
- Monitor system health dashboards
- Cannot modify any configuration
- Cannot access billing or sensitive data
- Typically assigned to support staff

**Auditor** - Compliance and security review:

- Full read access to audit logs
- Access to compliance reports
- Cannot modify anything
- Cannot view sensitive data (API keys, passwords)
- Access is logged for compliance
- Typically assigned to compliance officers or external auditors

### 1.3 Key Concepts

Concept	Description
<b>Tenant</b>	Organization with isolated data
<b>User</b>	End-user within a tenant
<b>Subscription</b>	Billing tier (1-7)
<b>Credits</b>	Currency for AI usage
<b>API Key</b>	Authentication for API access
<b>App</b>	Consumer application (Think Tank, etc.)

**Tenant Architecture Explained** A **Tenant** represents a complete organization using RADIANT. Each tenant has:

- **Complete Data Isolation:** All data is stored with tenant IDs and protected by PostgreSQL Row-Level Security (RLS). One tenant can never access another tenant's data, even if there's a bug in application code.
- **Separate Billing:** Each tenant has its own subscription, credit balance, and usage tracking. Costs are attributed to the correct tenant automatically.
- **Custom Configuration:** Tenants can customize model access, rate limits, and feature flags without affecting other tenants.
- **User Management:** Each tenant manages their own users, roles, and permissions independently.

**User vs Administrator** **Users** are end-users who interact with RADIANT-powered applications like Think Tank. They:

- Sign up and log in via Cognito
- Use AI models through the API or applications
- Have credits deducted for usage
- Cannot access the Admin Dashboard

**Administrators** manage the RADIANT platform itself. They:

- Access the Admin Dashboard
- Manage tenants, users, and billing
- Configure AI models and providers
- Have no credits (administrative access is separate)

**Credit System Explained** Credits are RADIANT's universal currency for AI usage:

- **1 credit = \$0.01 USD** (configurable per deployment)
- Different models cost different amounts based on their API pricing
- Credits are deducted in real-time as requests complete
- Tenants can purchase credits or receive them through subscriptions
- Credits can be tracked, audited, and reported on

**Example Credit Costs:** | Model | Cost per 1K tokens | |-----|-----| | GPT-4o | 5 credits input, 15 credits output | | GPT-4o-mini | 0.5 credits input, 1.5 credits output | | Claude 3.5 Sonnet | 3 credits input, 15 credits output | | Self-hosted Llama | 0.2 credits (all) |

**API Key Types** RADIANT supports multiple API key types:

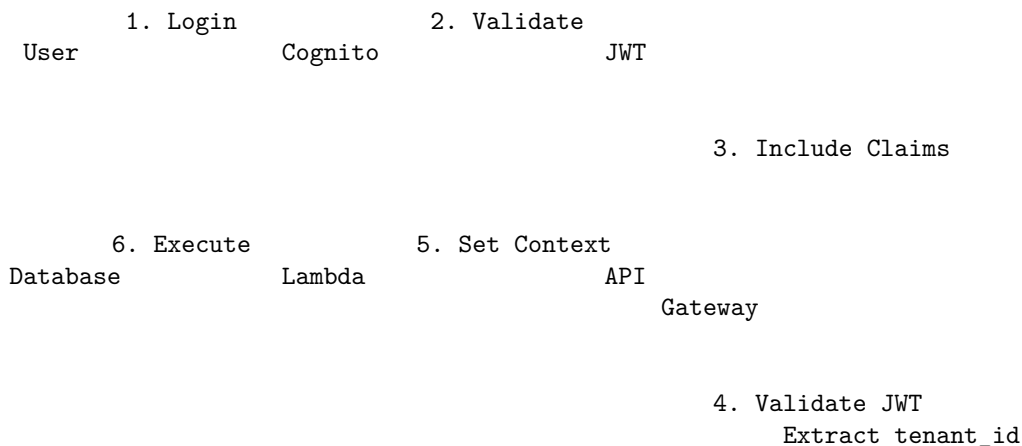
- **User API Keys:** Tied to a specific user, inherit user's permissions
- **Service API Keys:** For server-to-server communication, not tied to a user
- **Admin API Keys:** For administrative operations, require elevated permissions
- **Scoped Keys:** Limited to specific models, endpoints, or rate limits

## 1.4 Authentication Flow & Security

**Detailed Documentation:** For comprehensive authentication guides, see:

- Authentication Overview - Architecture and feature matrix
- Platform Admin Guide - Cognito management, global policies
- Security Architecture - Threat models, compliance
- Authentication API Reference - All auth endpoints

**JWT-Based Authentication** All authentication flows through Amazon Cognito with JWT tokens:



### Authentication Steps:

1. **Authentication:** Users authenticate via Amazon Cognito. Upon success, they receive a JSON Web Token (JWT).
2. **Token Validation:** The API Gateway Custom Authorizer validates the JWT signature using Cognito's public keys.
3. **Context Injection:** The Authorizer extracts the immutable `tenant_id` from the JWT claims and injects it into the request context.
4. **Lambda Context Setting:** The Lambda function sets PostgreSQL session variables before any database operation.
5. **RLS Enforcement:** All queries are automatically filtered by the database engine.

## JWT Claims Structure

```
{
  "sub": "user-uuid-here",
  "email": "user@example.com",
  "custom:tenant_id": "tenant-uuid-here",
  "custom:app_id": "thinktank",
  "custom:user_role": "member",
  "custom:tenant_role": "admin",
  "iat": 1735570000,
  "exp": 1735573600,
  "iss": "https://cognito-idp.us-east-1.amazonaws.com/us-east-1_XXXXX"
}
```

## Critical Security Invariant

**CRITICAL:** Application code **NEVER** accepts a `tenant_id` from the request body, query parameters, or custom headers. It **ONLY** trusts the ID injected by the API Gateway Authorizer from the cryptographically signed JWT.

This is the single most important security rule in the RADIANT architecture. Violation of this rule creates a critical vulnerability allowing tenant impersonation.

## Conditional Access Policies

Condition	Action
New device detected	Require MFA re-authentication
Unusual location	Challenge with security questions
After-hours access	Alert + additional logging
Multiple failed attempts	Temporary lockout + notification
API key from new IP	Rate limit + alert
Tenant status = suspended	Reject all requests
Tenant status = pending_deletion	Reject all requests

## 1.5 Compute Isolation (Lambda Tenant Isolation Mode)

RADIANT leverages **AWS Lambda Tenant Isolation Mode** (released 2025) to enforce strict boundaries between tenants sharing the same function code.

### How It Works

1. **Request Arrival** - Request arrives at Lambda with tenant identifier in event payload
2. **Tenant Detection** - Lambda service inspects the `tenant_id` field from the authorizer context
3. **Environment Binding** - Request routed to execution environment dedicated to that tenant
4. **Isolation Guarantee** - Warm environments are tenant-specific and never reused across tenants
5. **Resource Isolation** - `/tmp` directory and memory are cryptographically isolated per tenant

## Security Guarantees

Threat	Without Isolation Mode	With Isolation Mode
<code>/tmp</code> data leakage	Possible (shared container)	<b>Eliminated</b> (separate microVM)
Memory scraping	Possible (shared memory)	<b>Eliminated</b> (isolated memory)
Global variable pollution	Possible	<b>Eliminated</b>
Cold start timing attacks	Reduced	<b>Eliminated</b>

Threat	Without Isolation Mode	With Isolation Mode
Environment variable exposure	Possible	<b>Eliminated</b>

**Administrator Verification** To verify Tenant Isolation Mode is enabled:

1. AWS Console → Lambda → Function → Configuration → Concurrency
2. Check for "Tenant Isolation: Enabled" badge
3. Or via AWS CLI: `aws lambda get-function-configuration --function-name <name>`

## 1.6 Three-Layer Authentication Architecture (v5.1.1)

RADIANT implements a comprehensive three-layer authentication system for complete security coverage:

**Layer 1: End-User Authentication** End users authenticate via AWS Cognito User Pool with support for:

Feature	Description
<b>Email/Password</b>	Standard authentication with strong password policies
<b>MFA</b>	Optional or required TOTP, SMS, or email verification
<b>SSO Federation</b>	SAML 2.0 and OIDC integration for enterprise customers
<b>Domain Enforcement</b>	Force SSO for specific email domains

**Database Table: tenant\_users**

- RLS-protected tenant isolation
- Role-based access (standard\_user, tenant\_admin, tenant\_owner)
- Feature flags for app access (Think Tank, Curator, Tenant Admin)

**Layer 2: Platform Administrator Authentication** Platform admins use a separate Cognito Admin Pool with mandatory MFA:

Role	Permissions
<b>super_admin</b>	Full access, can delete tenants, manage other admins
<b>admin</b>	Tenant/user management, model configuration
<b>operator</b>	Read-only monitoring and support
<b>auditor</b>	Compliance logs and reports only

**Database Table: platform\_admins**

- No RLS (cross-tenant access required)
- Permission-based authorization
- Full audit trail of admin actions

**Layer 3: Service/Machine Authentication** API keys for server-to-server and automated integrations:

Feature	Description
<b>Scoped Access</b>	Fine-grained permissions (chat:read, embeddings:write, etc.)
<b>Rate Limiting</b>	Per-key and global rate limits
<b>IP Restrictions</b>	Whitelist allowed IP addresses
<b>Expiration</b>	Optional key expiration dates

Feature	Description
<b>Audit Logging</b>	All key usage tracked in <code>service_api_key_audit</code>

### API Key Scopes:

- `chat:read, chat:write` - Conversation access
- `models:read` - Model information
- `embeddings:write` - Vector embedding generation
- `files:read, files:write` - File upload/download
- `knowledge:read, knowledge:write` - Knowledge graph access
- `admin:read, admin:write` - Administrative operations

**Enterprise SSO Configuration** Navigate to **Settings** → **SSO Connections** to configure:

#### 1. SAML 2.0

- Upload IdP metadata XML or URL
- Configure attribute mapping
- Set up group-to-role mapping

#### 2. OIDC

- Enter issuer URL and client credentials
- Configure claim mapping
- Test connection before enabling

#### 3. Domain Enforcement

- Specify email domains requiring SSO
- Users from enforced domains cannot use password auth

### 1.7 Two-Factor Authentication (MFA) - v5.52.28

RADIANT enforces role-based Multi-Factor Authentication (MFA) using industry-standard TOTP (RFC 6238).

#### MFA Enforcement by Role

Role	MFA Required	Can Disable
<code>super_admin</code>	<b>Yes</b>	No
<code>admin</code>	<b>Yes</b>	No
<code>operator</code>	<b>Yes</b>	No
<code>auditor</code>	<b>Yes</b>	No
<code>tenant_admin</code>	<b>Yes</b>	No
<code>tenant_owner</code>	<b>Yes</b>	No
<code>standard_user</code>	No (future)	N/A

#### Critical Security Rules:

- Admin roles **cannot bypass** MFA enrollment
- Admin roles **cannot disable** MFA once enrolled
- MFA enrollment is enforced at login via full-screen gate

## MFA Enrollment Flow

1. **Login** - User enters email/password
2. **MFA Check** - System checks if role requires MFA
3. **Enrollment Gate** - If not enrolled, full-screen enrollment UI appears
4. **QR Code** - Scan with authenticator app (Google Authenticator, 1Password, Authy)
5. **Verification** - Enter 6-digit code to confirm
6. **Backup Codes** - Receive 10 one-time recovery codes
7. **Complete** - User can now access the dashboard

## Supported Authenticator Apps

- Google Authenticator
- Microsoft Authenticator
- 1Password
- Authy
- Any RFC 6238 compliant TOTP app

**Backup Codes** Each user receives 10 backup codes during enrollment:

- Format: XXXX-XXXX (8 alphanumeric characters)
- Each code can only be used **once**
- Warning displayed when < 3 codes remain
- Regenerate codes from **Settings** → **Security**

**Device Trust** Users can opt to "Trust this device for 30 days":

- Reduces MFA prompts on trusted devices
- Maximum 5 trusted devices per user
- Devices can be revoked from Settings
- Trust expires automatically after 30 days

## Lockout Policy

Condition	Action
3 failed MFA attempts	5-minute lockout
Lockout expires	Attempts reset
Successful verification	Attempts reset

**MFA Settings Page** Access MFA configuration at **Settings** → **Security**:

- **Status** - View MFA enabled/disabled status
- **Backup Codes** - View remaining count, regenerate codes
- **Trusted Devices** - List and revoke trusted devices
- **Audit Log** - View MFA-related events

## Database Tables

Table	Purpose
tenant_users.mfa_*	MFA columns for tenant users
platform_admins.mfa_*	MFA columns for platform admins
mfa_backup_codes	Hashed one-time recovery codes
mfa_trusted_devices	30-day device trust tokens
mfa_audit_log	Partitioned audit log for MFA events



Table	Purpose
-------	---------

## MFA API Endpoints

Endpoint	Method	Purpose
/api/v2/mfa/status	GET	Get MFA status
/api/v2/mfa/check	GET	Check if MFA required
/api/v2/mfa/enroll/start	POST	Start enrollment
/api/v2/mfa/enroll/verify	POST	Verify enrollment
/api/v2/mfa/verify	POST	Verify code during login
/api/v2/mfa/backup-codes/regenerate	POST	Regenerate backup codes
/api/v2/mfa/devices	GET	List trusted devices
/api/v2/mfa/devices/:id	DELETE	Revoke device

## 1.8 Internationalization & Multi-Language Search (v5.52.29)

RADIANT supports 18 languages for authentication UI and full-text search across all content.

### Supported Languages

Language	Code	Direction	Search Method
English	en	LTR	PostgreSQL <b>english</b>
Spanish	es	LTR	PostgreSQL <b>spanish</b>
French	fr	LTR	PostgreSQL <b>french</b>
German	de	LTR	PostgreSQL <b>german</b>
Portuguese	pt	LTR	PostgreSQL <b>portuguese</b>
Italian	it	LTR	PostgreSQL <b>italian</b>
Dutch	nl	LTR	PostgreSQL <b>dutch</b>
Polish	pl	LTR	PostgreSQL <b>simple</b>
Russian	ru	LTR	PostgreSQL <b>russian</b>
Turkish	tr	LTR	PostgreSQL <b>turkish</b>
Japanese	ja	LTR	<b>pg_bigm</b> bi-gram
Korean	ko	LTR	<b>pg_bigm</b> bi-gram
Chinese (Simplified)	zh-CN	LTR	<b>pg_bigm</b> bi-gram
Chinese (Traditional)	zh-TW	LTR	<b>pg_bigm</b> bi-gram
<b>Arabic</b>	ar	<b>RTL</b>	PostgreSQL <b>simple</b>
Hindi	hi	LTR	PostgreSQL <b>simple</b>
Thai	th	LTR	PostgreSQL <b>simple</b>
Vietnamese	vi	LTR	PostgreSQL <b>simple</b>

**CJK Search (Chinese, Japanese, Korean)** CJK languages require special handling because they don't use word boundaries:

- **pg\_bigm** extension provides bi-gram indexing
- Automatic language detection on content ingestion
- **search\_content()** function routes queries to appropriate search method
- GIN indexes on all searchable text columns

**RTL Support (Arabic)** Arabic users see proper right-to-left layouts:

- Automatic `dir="rtl"` on auth containers
- Flipped margins, paddings, and flex directions
- LTR preservation for emails, codes, passwords

## Database Migration 071

Column/Index	Table	Purpose
<code>detected_language</code>	<code>uds_conversations</code>	Auto-detected content language
<code>search_vector_simple</code>	<code>uds_conversations</code>	Fallback tsvector
<code>search_vector_english</code>	<code>uds_conversations</code>	Language-specific tsvector
<code>idx_*_bigm_*</code>	Multiple	GIN bi-gram indexes for CJK

## 1.9 System Overview Dashboard

Access real-time platform health at **System** → **Overview**:

### Component Health Monitoring

Component	Metrics Tracked
<b>LiteLLM Gateway</b>	Tasks, CPU, memory, requests/min
<b>Aurora PostgreSQL</b>	ACU, connections, IOPS
<b>ElastiCache Redis</b>	Memory, connections, cache hit rate
<b>Lambda Functions</b>	Invocations, concurrent executions, throttles
<b>API Gateway</b>	Requests/sec, error rates
<b>Cognito Pools</b>	Active users, sign-ins, failed logins

### Dashboard Features

- **Auto-Refresh:** Real-time updates every 30 seconds
- **Capacity Planning:** Utilization percentages with headroom alerts
- **Alert Management:** Active alerts with severity levels
- **Trend Analysis:** Up/down/stable indicators for all metrics

## 1.8 LiteLLM Gateway Configuration

Access gateway settings at **System** → **Gateway**:

### Auto-Scaling Parameters

Parameter	Description	Default
<b>Min Tasks</b>	Minimum running ECS tasks	2
<b>Max Tasks</b>	Maximum scaling limit	50
<b>Desired Tasks</b>	Initial task count	2
<b>Task CPU</b>	vCPU units per task (256=0.25 vCPU)	2048
<b>Task Memory</b>	Memory per task in MB	4096

### Scaling Thresholds

Threshold	Description	Default
<b>Target CPU Utilization</b>	Scale out when exceeded	70%
<b>Target Memory Utilization</b>	Scale out when exceeded	80%
<b>Target Requests/Target</b>	Requests per task before scaling	1000
<b>Scale Out Cooldown</b>	Wait between scale-out events	60s
<b>Scale In Cooldown</b>	Wait between scale-in events	300s

## Rate Limiting

Setting	Description	Default
<b>Global Rate Limit</b>	Max requests/second (all tenants)	10,000
<b>Per-Tenant Limit</b>	Max requests/minute per tenant	1,000
<b>Response Caching</b>	Cache identical requests	Enabled
<b>Cache TTL</b>	How long to cache responses	3600s

## Health Check Settings

Setting	Description	Default
<b>Health Check Path</b>	Endpoint for health checks	<code>/health</code>
<b>Check Interval</b>	Time between checks	30s
<b>Timeout</b>	Max wait for response	10s
<b>Unhealthy Threshold</b>	Failed checks before unhealthy	3

## 1.9 SSO Connections Management

Access SSO settings at **Settings** → **SSO**:

### Supported Protocols

Protocol	Use Case
<b>SAML 2.0</b>	Enterprise IdPs (Okta, Azure AD, OneLogin)
<b>OIDC</b>	Modern IdPs with OpenID Connect support

### Creating a SAML Connection

1. Navigate to **Settings** → **SSO**
2. Click **Add Connection**
3. Select **SAML 2.0** protocol
4. Enter:
  - **Connection Name**: Human-readable identifier
  - **Tenant ID**: UUID of the target tenant
  - **IdP Entity ID**: From your IdP metadata
  - **IdP SSO URL**: Login endpoint URL
  - **IdP Certificate**: X.509 certificate (PEM format)
5. Configure domain enforcement (optional)
6. Set default user role
7. Click **Create Connection**

## Creating an OIDC Connection

1. Navigate to **Settings** → **SSO**
2. Click **Add Connection**
3. Select **OIDC** protocol
4. Enter:
  - **Connection Name:** Human-readable identifier
  - **Tenant ID:** UUID of the target tenant
  - **Issuer URL:** Your IdP's issuer URL
  - **Client ID:** From your IdP
  - **Client Secret:** From your IdP
5. Configure domain enforcement (optional)
6. Set default user role
7. Click **Create Connection**

**Domain Enforcement** Force users with specific email domains to use SSO:

Enforced Domains: `example.com`, `corp.example.com`

Users with `@example.com` or `@corp.example.com` emails cannot use password authentication.

## Testing Connections

1. Find the connection in the list
  2. Click → **Test Connection**
  3. Verify the result:
    - **Success:** IdP is reachable and responding
    - **Failure:** Check configuration and IdP status
- 

## 2. Accessing the Admin Dashboard

### 2.1 URL and Login

1. Navigate to: `https://admin.your-domain.com`
2. Enter your email address
3. Enter your password
4. Complete MFA verification (required)

### 2.2 First Login

On first login:

1. You'll receive a temporary password via email
2. Enter the temporary password
3. Set a new password (12+ characters, mixed case, numbers, symbols)
4. Set up MFA using an authenticator app
5. You'll be redirected to the dashboard

### 2.3 Session Management

Setting	Value
<b>Session Duration</b>	8 hours
<b>Idle Timeout</b>	30 minutes
<b>Concurrent Sessions</b>	3 maximum
<b>Remember Device</b>	30 days

Setting	Value
---------	-------

2.4 Password Requirements

- Minimum 12 characters
- At least one uppercase letter
- At least one lowercase letter
- At least one number
- At least one special character
- Cannot reuse last 10 passwords

3. Dashboard Overview

3.1 Main Dashboard

The dashboard displays key metrics at a glance:

RADIANT Admin Dashboard

Welcome, Admin

Tenants	Users	Requests	Revenue
142	8,456	2.3M	\$45,230
+12%	+8%	+23%	+15%

Request Volume (7 days)

Mon

Tue

Wed

Thu

Fri

Sat

Sun

Recent Activity:

- New tenant: Acme Corp (2 minutes ago)
- Model enabled: claude-3-opus (15 minutes ago)
- Alert: High API error rate (1 hour ago)

3.2 Navigation Menu

Section	Description
<b>Dashboard</b>	Overview and metrics
<b>Tenants</b>	Tenant management
<b>Users</b>	User management
<b>Models</b>	AI model configuration
<b>Providers</b>	Provider management
<b>Billing</b>	Subscriptions and credits
<b>Storage</b>	Storage usage
<b>Orchestration</b>	Preference Engine settings
<b>Localization</b>	Translation management

Section	Description
<b>Configuration</b>	System settings
<b>Security</b>	Security monitoring
<b>Compliance</b>	Compliance reports
<b>Experiments</b>	A/B testing
<b>Cost</b>	Cost analytics
<b>Audit</b>	Audit logs
<b>Migrations</b>	Database migrations
<b>Notifications</b>	System alerts
<b>Settings</b>	Personal settings

## 4. Tenant Management

**New in v4.18.55:** Complete Tenant Management System (TMS) with lifecycle management, multi-tenant users, soft delete/restore, and compliance features.

### 4.1 Viewing Tenants

Navigate to **Tenants** to see all organizations:

Column	Description
<b>Name</b>	Organization display name and internal identifier
<b>Type</b>	Organization or Individual (phantom tenant)
<b>Tier</b>	Service tier (1-5: SEED, SPROUT, GROWTH, SCALE, ENTERPRISE)
<b>Status</b>	Active, Suspended, Pending, Pending Deletion, Deleted
<b>Compliance</b>	HIPAA, SOC2, GDPR badges if enabled
<b>Users</b>	Active user count with invited/suspended breakdown
<b>Created</b>	Creation date with relative time

### Tenant Types

Type	Description	Use Case
<b>Organization</b>	Company or team workspace	B2B customers
<b>Individual</b>	Personal workspace (phantom tenant)	B2C users, free tier

### Tenant Tiers

Tier	Name	Features
1	SEED	Basic features, standard support
2	SPROUT	Advanced features, email support
3	GROWTH	Premium features, dedicated KMS key
4	SCALE	Enterprise features, priority support
5	ENTERPRISE	Full features, custom SLA, dedicated resources

### 4.2 Creating a Tenant

1. Click ”+ **Create Tenant**” button

2. Fill in **Basic Information**:
  - **Internal Name**: URL-friendly identifier (e.g., `acme-corp`)
  - **Display Name**: Human-readable name (e.g., `Acme Corporation`)
  - **Type**: Organization or Individual
  - **Tier**: Select service tier (1-5)
3. Configure **Region & Compliance**:
  - **Primary Region**: `us-east-1`, `us-west-2`, `eu-west-1`, `eu-central-1`, `ap-southeast-1`
  - **Compliance Frameworks**: Check HIPAA, SOC2, GDPR as needed
  - **Domain** (optional): Custom domain for SSO
4. Set up **Initial Admin**:
  - **Admin Email**: Primary administrator email
  - **Admin Name**: Administrator display name
5. Click **"Create Tenant"**
- Note: HIPAA-enabled tenants automatically have a minimum 90-day retention period.

4.3 Tenant Statuses

Status	Description	Actions Available
Active	Normal operation	Edit, Suspend, Delete
Suspended	Temporarily disabled	Edit, Activate, Delete
Pending	Awaiting setup	Edit, Activate, Delete
Pending Deletion	Scheduled for deletion	Restore only
Deleted	Hard deleted	None (audit only)

4.3 Tenant Details

View comprehensive tenant information:

Tenant: Acme Corporation

OverviewUsersBillingSettings

Tenant ID:tn\_abc123xyz

Status:Active

Plan:Professional (Tier 4)

Created:2024-01-15

Last Active:2 minutes ago

Usage This Month:

API Requests:145,234

Tokens Used:12.5M

Storage:2.3 GB

Credits Used:\$1,234.56

[Edit]

[Suspend]

[Delete]

[Impersonate]

4.4 Soft Delete & Restore

Tenants are never immediately deleted. Instead, they go through a **retention period**:

## Deleting a Tenant

1. Navigate to **Tenants** → [Tenant]
2. Click the **Delete** button (or use dropdown menu)
3. Enter a **reason** for deletion (required)
4. Choose whether to **notify users**
5. Click **"Delete Tenant"**

The tenant will be marked as **"Pending Deletion"** with a scheduled hard-delete date.

## Retention Periods

Compliance	Minimum Retention	Default
None	7 days	30 days
GDPR	30 days	30 days
SOC2	30 days	30 days
HIPAA	90 days	90 days

## Restoring a Tenant

1. Find the tenant in the **Pending Deletions** section
2. Click **"Restore"**
3. Click **"Send Verification Code"**
4. Check your email for the 6-digit code
5. Enter the code and click **"Restore Tenant"**

**Security:** Verification codes expire after 15 minutes and allow only 3 attempts.

**Automatic Hard Delete** A scheduled job runs daily at **3:00 AM UTC** to:

- Process all tenants past their retention period
- Delete all S3 data for the tenant
- Schedule KMS key deletion (7-day pending window)
- Mark user records as deleted
- Send deletion confirmation emails

## 4.5 Multi-Tenant Users

Users can belong to **multiple tenants** with different roles:

### Membership Roles

Role	Permissions
<b>Owner</b>	Full access, can delete tenant, cannot be removed
<b>Admin</b>	Full access, can manage users
<b>Member</b>	Standard access to tenant resources
<b>Viewer</b>	Read-only access

## Adding Users to a Tenant

1. Navigate to **Tenants** → [Tenant] → **Users**
2. Click **" + Add User "**
3. Enter the user's email address
4. Select their role



5. Choose to send an invitation email (optional)
6. Click **"Add User"**

**No Orphan Users** The system **prevents orphan users** through a database trigger. When a user's last membership is removed, they are automatically soft-deleted. This ensures:

- Every active user has at least one tenant
- Billing is always attributed correctly
- No abandoned accounts consuming resources

## 4.6 Phantom Tenants

When a new user signs up without an invitation:

1. System creates an **Individual** type tenant automatically
2. User becomes the **Owner** of their phantom tenant
3. Tenant is named "[User's Name]'s Workspace"
4. User can later be invited to Organization tenants

This enables both **B2C** (individual users) and **B2B** (organizations) use cases.

## 4.7 Deletion Notifications

Users receive automatic emails when their tenant is scheduled for deletion:

Notification	Timing
<b>7-day warning</b>	7 days before deletion
<b>3-day warning</b>	3 days before deletion
<b>1-day warning</b>	1 day before deletion
<b>Deletion confirmation</b>	After hard delete

## 4.8 Compliance Features

**Risk Acceptances** For compliance frameworks, document risk acceptances:

1. Navigate to **Tenants** → [Tenant] → **Compliance**
2. View required controls for enabled frameworks
3. For controls that cannot be fully implemented:
  - Click **"Accept Risk"**
  - Enter risk description and mitigating controls
  - Provide business justification
  - Set expiration date
  - Submit for approval (if required)

**Compliance Reports** A scheduled job runs **monthly on the 1st at 4:00 AM UTC** generating:

- Tenant compliance breakdown (HIPAA, SOC2, GDPR)
- Risk acceptance status (pending, approved, expired)
- Retention compliance status
- Alerts for non-compliant tenants

## 4.9 Tenant Actions

Action	Description	Permission
<b>Edit</b>	Modify tenant settings	Admin

Action	Description	Permission
<b>Suspend</b>	Temporarily disable	Admin
<b>Delete</b>	Soft delete with retention	Admin
<b>Restore</b>	Restore from pending deletion	Admin + 2FA
<b>Manage Users</b>	Add/remove/modify memberships	Admin
<b>View Audit Log</b>	See tenant activity	Admin

#### 4.10 Data Isolation Architecture

RADIANT implements **six layers of tenant isolation** providing defense in depth:

##### TENANT ISOLATION LAYERS

###### LAYER 1: COMPUTE ISOLATION

AWS Lambda Tenant Isolation Mode (2025)

- Separate Firecracker microVM per tenant
- /tmp and memory NEVER shared between tenants
- Warm environments dedicated to specific tenant\_id

###### LAYER 2: NETWORK ISOLATION

API Gateway

- Per-tenant usage plans (rate limits by tier)
- AWS WAF rules (SQL injection, XSS, rate limiting)
- VPC endpoints for private connectivity (Enterprise tier)

###### LAYER 3: DATA ISOLATION (Polyglot Persistence)

Aurora PostgreSQL: Row-Level Security (RLS)

```
SET app.current_tenant_id = 'tenant-uuid';
```

```
POLICY: tenant_id = current_setting('app.current_tenant_id')::UUID
```

DynamoDB: IAM Leading Keys Condition

```
"dynamodb:LeadingKeys": ["TENANT#${tenant_id}#*"]
```

S3: Bucket Policy + Prefix Isolation

```
s3://bucket/tenants/${tenant_id}/
```

ElastiCache Redis: Key Prefix Namespacing

```
tenant:${tenant_id}:session:*
```

###### LAYER 4: ENCRYPTION ISOLATION

Tier 1-2: AWS-owned KMS key (shared, automatic rotation)

Tier 3: Customer Managed Key (CMK) per tenant

Tier 4-5: CMK with BYOK option (customer-controlled)

- Separate key for data at rest encryption
- Key deletion scheduled on tenant termination

## LAYER 5: IDENTITY ISOLATION

JWT Claims: tenant\_id extracted by API Gateway Authorizer  
NEVER trust tenant\_id from request body, headers, or query params  
Cognito User Pool with custom attributes per tenant  
Session tokens scoped to specific tenant context

## LAYER 6: CONTROL PLANE ISOLATION

Tenant Management Service (TMS) - Separate from Admin Dashboard

- Independent deployment and release cycle
- Higher security controls and MFA requirements
- Reduced blast radius from application compromise
- Internal-only API (not exposed to public internet)

### Storage Layer Isolation Summary

Data Type	Storage	Isolation Mechanism	Use Case
Users, Tenants, Conversations, Messages	Aurora PostgreSQL	Row-Level Security (RLS)	Complex queries, transactional
Sessions, Real-time presence	DynamoDB	Leading Keys + IAM	High velocity, TTL, global ta
Cache, Rate limiting	ElastiCache Redis	Key prefix namespacing	Sub-millisecond access
Media files, Exports	S3	Bucket policy + prefix	Large objects, CDN integrati
Audit logs	S3 + Athena	Object Lock + prefix	Immutable compliance record

### 4.11 Tenant Metadata Reference

Field	Type	Description
tenant_id	UUID	Primary key, system-generated (never user-provided)
name	VARCHAR(255)	Internal identifier (URL-friendly)
display_name	VARCHAR(255)	Human-readable name
type	ENUM	'organization' or 'individual'
status	ENUM	'active', 'suspended', 'pending', 'pending_deletion', 'deleted'
tier	INTEGER	Subscription tier (1-5): SEED, SPROUT, GROWTH, SCALE, ENTERPRISE
primary_region	VARCHAR(20)	AWS region for data residency (e.g., us-east-1)
compliance_mode	JSONB	Array: ['hipaa', 'soc2', 'gdpr']
retention_days	INTEGER	Custom retention period override (7-730 days)
deletion_scheduled_at	TIMESTAMPTZ	When hard delete will occur (if pending_deletion)
deletion_requested_by	UUID	User/admin who initiated deletion
stripe_customer_id	VARCHAR(255)	Billing integration identifier
kms_key_arn	VARCHAR(500)	Per-tenant encryption key (Tier 3+)
metadata	JSONB	Custom tenant metadata
settings	JSONB	Tenant-specific settings
created_at	TIMESTAMPTZ	Creation timestamp
updated_at	TIMESTAMPTZ	Last modification timestamp

#### 4.12 API Gateway Rate Limits by Tier

Tier	Requests/Second	Burst	Monthly Quota	Price Point
Tier 1 (SEED)	10	50	100,000	\$200/month
Tier 2 (SPROUT)	50	200	1,000,000	\$1,000/month
Tier 3 (GROWTH)	200	500	10,000,000	\$5,000/month
Tier 4 (SCALE)	1,000	2,000	100,000,000	\$25,000/month
Tier 5 (ENTERPRISE)	Custom	Custom	Unlimited	\$150,000+/month

**Administrator Action:** Usage plan assignment happens automatically based on tenant tier. Override via Admin Dashboard → Tenant → API Settings.

#### 4.13 TMS API Endpoints Reference

##### Tenant CRUD

Method	Path	Description	Auth
POST	/internal/tms/tenants	Create tenant	Internal Service Token
GET	/internal/tms/tenants/{id}	Get tenant	Internal Service Token
PUT	/internal/tms/tenants/{id}	Update tenant	Internal Service Token
DELETE	/internal/tms/tenants/{id}	Soft delete tenant	Internal Service Token
POST	/internal/tms/tenants/{id}/restore	Restore tenant	Internal Service Token + 2FA
POST	/internal/tms/phantom-tenant	Create phantom tenant	Internal Service Token

##### User Membership

Method	Path	Description	Auth
GET	/internal/tms/tenants/{id}/users	List tenant users	Internal Service Token
POST	/internal/tms/tenants/{id}/users	Add user to tenant	Internal Service Token
PUT	/internal/tms/tenants/{id}/users/{userId}	Update membership	Internal Service Token
DELETE	/internal/tms/tenants/{id}/users/{userId}	Remove user	Internal Service Token

##### Request/Response Examples Create Tenant Request:

```
{
  "name": "acme-corp",
  "displayName": "Acme Corporation",
  "type": "organization",
  "tier": 3,
  "primaryRegion": "us-east-1",
  "complianceMode": ["hipaa", "soc2"],
  "adminEmail": "admin@acme.com",
  "adminName": "John Smith"
}
```

##### Create Tenant Response:

```
{
  "success": true,
  "data": {
    "tenant": {
```

```

    "id": "550e8400-e29b-41d4-a716-446655440000",
    "name": "acme-corp",
    "displayName": "Acme Corporation",
    "type": "organization",
    "status": "active",
    "tier": 3,
    "kmsKeyArn": "arn:aws:kms:us-east-1:123456789012:key/..."
  },
  "adminUser": {
    "id": "user-uuid-here",
    "email": "admin@acme.com"
  },
  "membership": {
    "role": "owner",
    "status": "active"
  }
}
}
}

```

#### Soft Delete Request:

```

{
  "initiatedBy": "admin-uuid",
  "reason": "Customer requested account closure",
  "notifyUsers": true
}

```

#### Soft Delete Response:

```

{
  "success": true,
  "data": {
    "tenantId": "550e8400-e29b-41d4-a716-446655440000",
    "status": "pending_deletion",
    "deletionScheduledAt": "2025-01-30T03:00:00Z",
    "retentionDays": 30,
    "affectedUsers": {
      "total": 47,
      "willBeDeleted": 42,
      "willRemain": 5
    },
    "notificationsSent": true
  }
}

```

#### 4.14 Scheduled Jobs

Job	Schedule	Description
<b>Hard Delete</b>	Daily 3:00 AM UTC	Processes tenants past retention period
<b>Deletion Notifications</b>	Daily 9:00 AM UTC	Sends 7/3/1 day warnings
<b>Orphan Check</b>	Weekly Sunday 2:00 AM UTC	Safety net for orphan user cleanup
<b>Compliance Report</b>	Monthly 1st 4:00 AM UTC	Generates compliance status report

#### Hard Delete Job Details Process:

1. Query tenants where `status = 'pending_deletion'` AND `deletion_scheduled_at < NOW()`
2. For each tenant (up to 10 per run to avoid Lambda timeout):
  - Count users that will be deleted vs retained (multi-tenant users)
  - Delete all tenant memberships (triggers orphan user check)
  - Delete S3 data under `tenants/{tenant_id}/` prefix
  - Schedule KMS key for deletion (7-30 day AWS delay)
  - Update tenant status to `'deleted'`
  - Send deletion confirmation emails
  - Log to immutable audit trail

**Administrator Note:** Hard deletes cannot be reversed. Monitor the audit log for any unexpected deletions.

#### 4.15 Database Schema (Migration 126)

##### tenant\_user\_memberships Table

```
CREATE TABLE tenant_user_memberships (
  id UUID PRIMARY KEY DEFAULT gen_random_uuid(),
  tenant_id UUID NOT NULL REFERENCES tenants(id) ON DELETE CASCADE,
  user_id UUID NOT NULL REFERENCES users(id) ON DELETE CASCADE,
  role VARCHAR(50) DEFAULT 'member'
    CHECK (role IN ('owner', 'admin', 'member', 'viewer')),
  status VARCHAR(20) DEFAULT 'active'
    CHECK (status IN ('active', 'suspended', 'invited')),
  joined_at TIMESTAMPTZ DEFAULT NOW(),
  invited_by UUID REFERENCES users(id),
  invitation_token VARCHAR(255),
  invitation_expires_at TIMESTAMPTZ,
  last_active_at TIMESTAMPTZ,
  created_at TIMESTAMPTZ DEFAULT NOW(),
  updated_at TIMESTAMPTZ DEFAULT NOW(),
  UNIQUE(tenant_id, user_id)
);
```

##### tms\_audit\_log Table

```
CREATE TABLE tms_audit_log (
  id UUID PRIMARY KEY DEFAULT gen_random_uuid(),
  tenant_id UUID REFERENCES tenants(id) ON DELETE SET NULL,
  user_id UUID REFERENCES users(id) ON DELETE SET NULL,
  admin_id UUID REFERENCES administrators(id) ON DELETE SET NULL,
  action VARCHAR(100) NOT NULL,
  resource_type VARCHAR(50),
  resource_id VARCHAR(255),
  old_value JSONB,
  new_value JSONB,
  ip_address INET,
  user_agent TEXT,
  trace_id VARCHAR(100),
  created_at TIMESTAMPTZ DEFAULT NOW()
);
```

##### tms\_verification\_codes Table

```
CREATE TABLE tms_verification_codes (
  id UUID PRIMARY KEY DEFAULT gen_random_uuid(),
```

```

    user_id UUID REFERENCES users(id) ON DELETE CASCADE,
    admin_id UUID REFERENCES administrators(id) ON DELETE CASCADE,
    operation VARCHAR(50) NOT NULL
        CHECK (operation IN ('restore_tenant', 'hard_delete', 'transfer_ownership', 'compliance_override')),
    resource_id UUID NOT NULL,
    code_hash VARCHAR(64) NOT NULL,
    attempts INTEGER DEFAULT 0,
    max_attempts INTEGER DEFAULT 3,
    expires_at TIMESTAMPTZ NOT NULL,
    verified_at TIMESTAMPTZ,
    created_at TIMESTAMPTZ DEFAULT NOW(),
    CONSTRAINT chk_user_or_admin CHECK (user_id IS NOT NULL OR admin_id IS NOT NULL)
);

```

### Orphan Prevention Trigger

```

CREATE OR REPLACE FUNCTION tms_prevent_orphan_users()
RETURNS TRIGGER AS $$
DECLARE
    remaining_memberships INTEGER;
BEGIN
    SELECT COUNT(*) INTO remaining_memberships
    FROM tenant_user_memberships
    WHERE user_id = OLD.user_id
    AND id != OLD.id
    AND status != 'invited';

    IF remaining_memberships = 0 THEN
        UPDATE users SET status = 'deleted', updated_at = NOW()
        WHERE id = OLD.user_id;

        INSERT INTO tms_audit_log (tenant_id, user_id, action, resource_type, resource_id)
        VALUES (OLD.tenant_id, OLD.user_id, 'user_auto_deleted_orphan', 'user', OLD.user_id::text);
    END IF;

    RETURN OLD;
END;
$$ LANGUAGE plpgsql;

CREATE TRIGGER tms_check_orphan_on_membership_delete
AFTER DELETE ON tenant_user_memberships
FOR EACH ROW EXECUTE FUNCTION tms_prevent_orphan_users();

```

## 4.16 Troubleshooting

Issue	Cause	Resolution
Cannot delete tenant	Tenant has HIPAA compliance	Verify 90-day minimum retention is met
Restore code not received	Email delivery failed	Check spam folder, verify email address, resend code
Restore code invalid	Code expired or max attempts	Request new code (3 attempt limit, 15 min expiry)
User not deleted with tenant	User has other tenant memberships	Expected behavior - user retained for other tenants
Phantom tenant not created	User email already exists	User will be associated with existing account
KMS key not created	Tenant tier below 3	Upgrade to Tier 3+ for dedicated KMS key
Cannot remove last owner	Business rule enforcement	Transfer ownership first or delete tenant

---

## 5. User & Administrator Management

### 5.1 Administrator Roles

Role	Dashboard Access	API Access	Billing	Audit
<b>Super Admin</b>	Full	Full	Full	Full
<b>Admin</b>	Full	Full	Read	Read
<b>Operator</b>	Read	Read	None	Read
<b>Auditor</b>	Logs only	None	None	Full

### 5.2 Managing Administrators

Navigate to **Administrators** to:

1. **Invite New Admin:**

- Click **" + Invite Administrator "**
- Enter email address
- Select role
- Click **"Send Invitation"**

2. **Modify Admin:**

- Click on administrator row
- Edit role or permissions
- Click **"Save Changes"**

3. **Remove Admin:**

- Click **"Remove"** button
- Confirm removal
- Admin's sessions are invalidated immediately

### 5.3 Viewing Tenant Users

Navigate to **Tenants** → **[Tenant]** → **Users** to see:

Field	Description
<b>Email</b>	User email
<b>Name</b>	Display name
<b>Role</b>	Tenant role
<b>Status</b>	Active/Invited/Disabled
<b>Last Login</b>	Last authentication
<b>API Keys</b>	Number of active keys

### 5.4 User Actions

Action	Description
<b>Reset Password</b>	Send password reset email
<b>Disable</b>	Prevent login
<b>Enable</b>	Restore access
<b>Delete</b>	Remove user data



Action	Description
<b>View Sessions</b>	See active sessions

## 6. AI Model Configuration

### 6.1 Model Registry

Navigate to **Models** to see all available models:

AI Models 106 Total

Filter: [All ] Category: [All ] Status: [Enabled ]

Model	Provider	Category	Tier	Status
gpt-4o	OpenAI	Chat	1	Enabled
gpt-4o-mini	OpenAI	Chat	1	Enabled
claude-3-opus	Anthropic	Chat	2	Enabled
claude-3-sonnet	Anthropic	Chat	1	Enabled
gemini-pro	Google	Chat	1	Enabled
llama-3.1-70b	Self-Host	Chat	3	Enabled
whisper-large	Self-Host	Audio	3	Disabled

[+ Add Model] [Import Models] [Export Config]

### 6.2 Model Categories

Category	Description	Example Models
<b>Chat/LLM</b>	Text generation	GPT-4o, Claude 3, Gemini
<b>Embedding</b>	Vector embeddings	text-embedding-3-large
<b>Vision</b>	Image understanding	GPT-4V, Claude Vision
<b>Audio</b>	Speech-to-text	Whisper, Deepgram
<b>Image</b>	Image generation	DALL-E 3, Stable Diffusion
<b>Code</b>	Code generation	Codestral, DeepSeek Coder
<b>Scientific</b>	Research models	BioGPT, ChemLLM

**Category Details Chat/LLM (Large Language Models):** The core of RADIANT. These models handle conversational AI, content generation, summarization, and general-purpose text tasks. They're the most commonly used and include flagship models from OpenAI, Anthropic, Google, and open-source alternatives.

**Embedding Models:** Convert text into numerical vectors for semantic search, similarity matching, and retrieval-augmented generation (RAG). Essential for building knowledge bases and search functionality. Vectors are typically 1536-3072 dimensions.

**Vision Models:** Analyze images, extract text (OCR), describe visual content, and answer questions about images. Increasingly important for document processing, accessibility, and multimodal applications.

**Audio Models:** Transcribe speech to text, translate audio, and identify speakers. Whisper is the most popular, offering excellent accuracy across 99 languages. Used for meeting transcription, accessibility, and voice interfaces.

**Image Generation:** Create images from text descriptions. DALL-E 3 offers the best prompt following, while Stable Diffusion provides more customization options. Consider content policies when enabling these.

**Code Models:** Specialized for programming tasks including code generation, explanation, debugging, and refactoring. Some are fine-tuned on specific languages or frameworks.

**Scientific Models:** Domain-specific models trained on scientific literature. Useful for research applications but require careful evaluation for accuracy.

### 6.3 Model Configuration

Click on a model to configure:

Setting	Description
<b>Enabled</b>	Available for use
<b>Min Tier</b>	Minimum subscription tier
<b>Rate Limits</b>	Requests per minute
<b>Max Tokens</b>	Maximum context/output
<b>Temperature Range</b>	Allowed temperature values
<b>Price Override</b>	Custom pricing

**Configuration Settings Explained** **Enabled:** When disabled, the model is hidden from users and API requests return "model not found". Use this to temporarily remove models during maintenance or to restrict access to specific models.

**Min Tier:** Sets the minimum subscription tier required to access this model. For example, setting GPT-4 to Tier 2 means Free tier users cannot use it. This helps control costs and create upgrade incentives.

**Rate Limits:** Controls requests per minute per user for this model. Prevents abuse and ensures fair access. Set based on the provider's rate limits and your capacity:

- Conservative: 10-20 requests/minute
- Standard: 50-100 requests/minute
- High: 200+ requests/minute (requires provider rate limit increases)

**Max Tokens:** Limits context window and output length. Useful for controlling costs since longer contexts cost more. Set based on use case:

- Short tasks (Q&A): 4,096 tokens
- Medium tasks (writing): 16,384 tokens
- Long tasks (analysis): 32,768+ tokens

**Temperature Range:** Restricts the temperature parameter users can set. Temperature controls randomness:

- 0.0: Deterministic, consistent outputs
- 0.7: Balanced creativity and consistency
- 1.0+: More creative, less predictable

Restricting range (e.g., 0.0-1.0) prevents users from setting extreme values that produce poor results.

**Price Override:** Allows custom pricing different from the default. Useful for:

- Offering discounts on specific models
- Increasing prices for premium models

- Matching competitor pricing
- A/B testing pricing strategies

## 6.4 Self-Hosted Models

For Tier 3+ deployments:

1. Navigate to **Models** → **Self-Hosted**
2. Click "+ Add Self-Hosted Model"
3. Configure:
  - **Model ID**: Unique identifier
  - **SageMaker Endpoint**: Endpoint name
  - **Instance Type**: ml.g5.xlarge, etc.
  - **Auto-Scaling**: Min/max instances
4. Deploy model to SageMaker

## 6.5 Thermal States (Self-Hosted)

State	Description	Response Time
<b>HOT</b>	Always running	<100ms
<b>WARM</b>	Scaled down	<5s
<b>COLD</b>	Stopped	30-60s
<b>OFF</b>	Disabled	N/A

# 7. Provider Management

## 7.1 External Providers

Navigate to **Providers** to manage API integrations:

Provider	Models	Status	Health
<b>OpenAI</b>	12	Configured	99.9%
<b>Anthropic</b>	6	Configured	99.8%
<b>Google AI</b>	8	Configured	99.7%
<b>xAI</b>	2	Configured	99.5%
<b>DeepSeek</b>	4	Not configured	-

## 7.2 Adding Provider Credentials

1. Click on provider name
2. Click "Configure"
3. Enter API credentials:
  - **API Key**: Provider API key
  - **Organization ID**: (if applicable)
  - **Base URL**: (for custom endpoints)
4. Click "Test Connection"
5. Click "Save"

## 7.3 Provider Health Monitoring

View real-time provider health:

Provider Health: OpenAI

Status: Healthy  
Uptime (30d): 99.94%  
Avg Latency: 245ms  
P95 Latency: 520ms  
Error Rate: 0.02%

Last 24 Hours:

12am                  6am                  12pm                  6pm                  12am

7.4 Fallback Configuration

Configure provider fallbacks:

- 1. Navigate to **Providers** → **Fallbacks**
- 2. Set priority order for each model category
- 3. Configure automatic failover rules
- 4. Set retry policies

8. Billing & Subscriptions

8.1 Subscription Tiers

Tier	Name	Monthly	Features
1	Free	\$0	Basic models, 1K requests
2	Starter	\$29	More models, 10K requests
3	Professional	\$99	All external models, 100K requests
4	Business	\$299	Priority support, 500K requests
5	Enterprise	\$999	Self-hosted, unlimited
6	Enterprise+	Custom	Custom SLAs, dedicated support
7	Ultimate	Custom	On-premise options

8.2 Credit System

Credits are the universal currency for AI usage:

Model Type	Cost per 1M Tokens
GPT-4o	500 credits
GPT-4o-mini	50 credits
Claude 3 Opus	600 credits
Claude 3 Sonnet	150 credits
Self-hosted	20 credits

8.3 Managing Subscriptions

Navigate to **Billing** → **Subscriptions**:

1. View current subscription
2. Upgrade/downgrade tier
3. Add credit packages
4. View invoices
5. Update payment method

## 8.4 Usage Reports

Generate usage reports:

1. Navigate to **Billing** → **Reports**
2. Select date range
3. Choose grouping (by tenant/model/user)
4. Export as CSV/PDF

## 8.5 Billing Alerts

Configure alerts for:

- Credit balance low
  - Usage spike
  - Approaching quota
  - Failed payments
- 

# 9. Storage Management

## 9.1 Storage Overview

Navigate to **Storage** to monitor:

### Storage Overview

Total Used:            234.5 GB of 500 GB (47%)

#### By Type:

Documents:	120.3 GB (51%)
Images:	45.2 GB (19%)
Audio:	38.7 GB (17%)
Video:	22.1 GB (9%)
Other:	8.2 GB (4%)

#### Top Tenants:

1. Acme Corp	45.2 GB
2. TechStart	32.1 GB
3. DataCo	28.4 GB

## 9.2 Storage Tiers

Tier	Included	Additional
Free	1 GB	N/A
Starter	10 GB	\$0.10/GB
Professional	100 GB	\$0.08/GB
Business	500 GB	\$0.05/GB
Enterprise	2 TB	\$0.03/GB

### 9.3 File Management

Manage uploaded files:

- View file metadata
- Download files
- Delete files
- Set retention policies

## 10. Orchestration & Preference Engine

### 10.1 Cognitive Router (formerly Brain Router)

The Cognitive Router automatically selects optimal models:

Factor	Weight	Description
<b>Cost</b>	30%	Price optimization
<b>Quality</b>	30%	Output quality
<b>Speed</b>	20%	Response latency
<b>Availability</b>	20%	Provider health

### 10.2 Neural Patterns

Configure orchestration patterns:

Pattern	Description	Use Case
<b>Single</b>	One model	Simple requests
<b>Fallback</b>	Primary + backup	High availability
<b>Parallel</b>	Multiple simultaneous	Consensus
<b>Chain</b>	Sequential models	Complex tasks

### 10.3 Workflow Templates

Create reusable workflows:

1. Navigate to **Orchestration** → **Workflows**
2. Click "+ New Workflow"
3. Define steps and conditions
4. Set triggers and parameters
5. Save and activate

### 10.4 Orchestration Methods

**Location:** Admin Dashboard → Orchestration → Methods

Manage the 70+ built-in orchestration methods that power workflow steps.

## Method Categories

Category	Count	Examples
<b>Generation</b>	3	Chain-of-Thought, Iterative Refinement
<b>Evaluation</b>	6	Multi-Judge Panel, G-Eval, Pairwise
<b>Synthesis</b>	5	Mixture of Agents, LLM-Blender
<b>Verification</b>	8	SelfCheckGPT, CiteFix, PRM
<b>Debate</b>	5	Sparse Debate, HAH-Delphi
<b>Aggregation</b>	4	Self-Consistency, GEDI Electoral
<b>Routing</b>	7	RouteLLM, FrugalGPT, Pareto, C3PO, AutoMix
<b>Uncertainty</b>	6	Semantic Entropy, SE Probes, Kernel Entropy
<b>Hallucination</b>	3	Multi-Method Detection, MetaQA
<b>Human-in-Loop</b>	3	HITL Review, Active Learning

## System vs User Methods

Type	Editable	Description
<b>System</b>	Parameters only	Built-in methods with locked definitions
<b>User</b>	Full (future)	Custom tenant methods (planned feature)

System methods show a "System" badge in the UI. Only `defaultParameters` and `isEnabled` can be modified.

## Setting Default Parameters

1. Navigate to **Orchestration** → **Methods**
2. Select a method from the list
3. Adjust parameters in the right panel
4. Click **Save** to update defaults

**Example Parameters** (vary by method):

Parameter	Type	Used By	Description
<code>sample_count</code>	integer	Entropy, Consistency	Number of response samples
<code>threshold</code>	number	Routing, Detection	Decision threshold
<code>temperature</code>	number	Generation, Debate	Sampling temperature
<code>confidence_threshold</code>	number	Cascade, HITL	Escalation trigger
<code>budget_cents</code>	number	Pareto Routing	Cost constraint
<code>kernel</code>	enum	Kernel Entropy	RBF, linear, polynomial
<code>topology</code>	enum	Sparse Debate	ring, star, tree, full

See `THINKTANK-ADMIN-GUIDE.md` Section 34.5 for complete parameter reference by method.

## Parameter Inheritance

Admin Default → Workflow Step → User Template Override

↓                      ↓                      ↓  
Tenant-wide      Per-workflow      Per-user template

Parameters set at the Admin level apply to all workflows. User workflow templates can override these per-template.

## 10.5 Orchestration Security (RLS)

All orchestration tables are protected by PostgreSQL Row-Level Security (RLS) with tenant isolation.

### Security Model Summary

Table	Read Access	Write Access
orchestration_methods	All authenticated users	Super admin only
orchestration_workflows	System: all users; User-created: own tenant only	Own tenant only
workflow_method_bindings	Follows parent workflow	Follows parent workflow
workflow_customizations	Own tenant only	Own tenant only
orchestration_executions	Own tenant (+ super admin)	Own tenant only
orchestration_step_executions	Follows parent execution	Follows parent execution
user_workflow_templates	Own + shared + approved public	Own only (tenant admin can manage)

**Session Variables** The API sets these PostgreSQL session variables before any database operation:

Variable	Purpose	Set By
app.current_tenant_id	Current tenant UUID	JWT claims
app.current_user_id	Current user UUID	JWT claims
app.is_tenant_admin	Tenant admin flag	JWT claims
app.is_super_admin	Super admin flag	JWT claims
app.client_ip	Client IP address	Request context
app.user_agent	Browser user agent	Request headers

**Helper Functions** Three helper functions are available for workflow access control:

```
-- Check if current user can view a workflow
SELECT can_access_workflow('workflow-uuid');

-- Check if current user can modify a workflow
SELECT can_modify_workflow('workflow-uuid');

-- Get all accessible workflows with permissions
SELECT * FROM get_accessible_workflows();
```

**Audit Trail** All modifications to workflows are logged to `orchestration_audit_log`:

Column	Description
table_name	Which table was modified
record_id	UUID of the modified record
action	INSERT, UPDATE, or DELETE
old_data	Previous state (JSONB)
new_data	New state (JSONB)
tenant_id	Tenant making the change
user_id	User making the change
ip_address	Client IP address
user_agent	Browser user agent
changed_at	Timestamp of change



**User Workflow Template Sharing** Users can share their workflow templates at three levels:

1. **Private** (default): Only the creator can see/use the template
2. **Shared within Tenant** (`is_shared = true`): All users in the tenant can see/use
3. **Public** (`is_public = true`): Visible to all tenants after super admin approval

To approve a public template (super admin only):

```
UPDATE user_workflow_templates
SET share_approved_at = NOW(), share_approved_by = 'admin-uuid'
WHERE template_id = 'template-uuid' AND is_public = true;
```

Migration: V2026\_01\_10\_003\_\_orchestration\_qls\_security.sql

---

## 11. Pre-Prompt Learning

**Location:** Admin Dashboard → Orchestration → Pre-Prompts

The Pre-Prompt Learning system tracks and learns from the effectiveness of pre-prompts (system prompts) used by the AGI Brain. It uses **attribution analysis** to determine what factor actually caused issues - not just the pre-prompt.

### 11.1 Overview Dashboard

The dashboard shows key metrics:

Metric	Description
<b>Templates</b>	Active/total pre-prompt templates
<b>Total Uses</b>	Pre-prompt instances used
<b>Avg Rating</b>	Average user rating (1-5)
<b>Thumbs Up Rate</b>	Percentage of positive feedback
<b>Exploration Rate</b>	Rate of trying non-optimal templates to learn

### 11.2 Attribution Analysis

When users report issues, the system analyzes which factor was responsible:

Factor	When Blamed
<b>Pre-prompt</b>	System instructions wrong
<b>Model</b>	Model selection inappropriate
<b>Mode</b>	Orchestration mode wrong
<b>Workflow</b>	Pattern did not fit task
<b>Domain</b>	Domain detection incorrect
<b>Other</b>	External factors

The attribution pie chart shows the distribution of blame across factors.

### 11.3 Template Management

Navigate to the **Templates** tab to:

- View all pre-prompt templates
- See usage statistics and success rates
- Check which orchestration modes each template supports
- View average feedback scores

## 11.4 Adjusting Weights

Click the **sliders icon** on any template to adjust weights:

Weight	Default	Description
Base Effectiveness	0.5	Starting score before bonuses
Domain Weight	0.2	Bonus for matching domain
Mode Weight	0.2	Bonus for matching mode
Model Weight	0.2	Bonus for compatible model
Complexity Weight	0.15	Bonus for complexity match
Task Type Weight	0.15	Bonus for task type match
Feedback Weight	0.1	Historical feedback influence

### Score Formula:

Final Score = Base + (Domain × DomainWeight) + (Mode × ModeWeight) +  
(Model × ModelWeight) + (Complexity × ComplexityWeight) +  
(TaskType × TaskTypeWeight) + FeedbackAdjustment

## 11.5 Learning Configuration

Configure learning behavior:

- **Learning Enabled:** Toggle on/off
- **Exploration Rate:** Percentage of requests using non-optimal templates to gather data (default: 10%)
- **Exploration Decay:** How quickly exploration decreases (default: 0.99)
- **Minimum Exploration:** Floor for exploration rate (default: 1%)

## 11.6 Recent Feedback

The **Feedback** tab shows recent user feedback with:

- Rating (1-5 stars)
- Attribution label (what got blamed)
- User comments
- Timestamp

See Pre-Prompt Learning System Documentation for full details.

---

## 12. Localization

### 12.1 Translation Registry

**Location:** Admin Dashboard → Localization → Translation Registry

The Translation Registry provides comprehensive management of all UI strings across RADIANT applications with tenant-specific override capabilities.

#### Key Features:

- **100+ Pre-Seeded Strings:** Common UI elements, errors, dialogs, validation messages
- **App-Scoped Categories:** radiant\_admin, thinktank\_admin, thinktank, curator, common
- **Tenant Overrides:** Custom text per tenant with protection from auto-updates
- **18 Language Support:** Full coverage across all supported languages

## 12.2 Supported Languages

Language	Code	Flag	RTL
English	en		No
Spanish	es		No
French	fr		No
German	de		No
Portuguese	pt		No
Italian	it		No
Dutch	nl		No
Polish	pl		No
Russian	ru		No
Turkish	tr		No
Japanese	ja		No
Korean	ko		No
Chinese (Simplified)	zh-CN		No
Chinese (Traditional)	zh-TW		No
Arabic	ar		<b>Yes</b>
Hindi	hi		No
Thai	th		No
Vietnamese	vi		No

## 12.3 Tenant Translation Overrides

Tenant administrators can override any system string with custom text:

### Override Workflow:

1. Navigate to **Localization** → **Translation Registry**
2. Select target language from dropdown
3. Browse or search for strings to customize
4. Click **Edit** to create an override
5. Enter custom text and save

### Protection System:

- **Protected Overrides** (default): Won't be updated by automatic translation
- **Unprotected Overrides**: May be updated when system translations improve
- **Revert**: Delete override to return to system translation

**Admin UI Tabs:** | Tab | Purpose | |----|-----| | **Registry** | Browse all strings, create overrides | | **Your Overrides** | Manage existing overrides, toggle protection | | **Coverage** | View translation coverage by language |

## 12.4 Translation Override API

Base URL: /api/admin/localization

### List Registry Entries

GET /api/admin/localization/registry?app\_id=thinktank&category=chat&search=error&page=1&limit=50

### Get Entry with All Translations

GET /api/admin/localization/registry/:id

### List Tenant Overrides

GET /api/admin/localization/overrides?language\_code=es&protected=true

### Create/Update Override

POST /api/admin/localization/overrides

Content-Type: application/json

```
{
  "registry_id": 42,
  "language_code": "es",
  "override_text": "Texto personalizado",
  "is_protected": true
}
```

### Delete Override (Revert to System)

DELETE /api/admin/localization/overrides/:id

### Toggle Protection

PATCH /api/admin/localization/overrides/:id/protection

Content-Type: application/json

```
{
  "is_protected": false
}
```

### Get Translation Bundle (with Overrides Applied)

GET /api/admin/localization/bundle/es?app\_id=thinktank

Response:

```
{
  "languageCode": "es",
  "translations": {
    "thinktank.chat.placeholder": "Escribe tu mensaje...",
    "common.buttons.save": "Guardar"
  },
  "overrideCount": 5,
  "overrideKeys": ["thinktank.chat.placeholder"]
}
```

## 12.5 Tenant Localization Config

Configure language settings per tenant:

GET /api/admin/localization/config

PUT /api/admin/localization/config

**Configuration Options:** | Setting | Default | Description | |-----|-----|-----| | default\_language | en | Default language for new users | | enabled\_languages | ['en'] | Languages available to users | | allow\_user\_language\_selection | true | Let users choose their language | | enable\_ai\_translation | true | Use AI for missing translations | | brand\_name | null | Custom brand name for strings |

## 12.6 Database Schema

Table	Purpose
localization_registry	Source strings with keys, context, category
localization_translations	System translations per language
tenant_translation_overrides	Tenant-specific overrides with protection
tenant_localization_config	Per-tenant language configuration
translation_audit_log	Change history for compliance

## 12.7 String Categories

Category	Examples
buttons	Save, Cancel, Delete, Edit, Submit
errors	Network error, Session expired, Validation
dialogs	Confirm, Delete confirmation, Unsaved changes
validation	Required field, Invalid email, Min/max length
toasts	Success messages, Error messages
time	Relative time (minutes ago, hours ago)
pagination	Showing X to Y, Previous, Next
tables	No data, Loading, Search
upload	Drag and drop, Max size, Allowed types
auth	Sign in, Sign out, Forgot password
navigation	Dashboard, Settings, Users
chat	Placeholder, Send, Thinking

## 12.8 Think Tank Admin Localization

Tenant administrators access a simplified localization UI at **Administration** → **Localization**:

- **Your Overrides**: View and manage custom translations
- **Browse Strings**: Search Think Tank-specific strings
- **Configuration**: Set default language and enabled languages

## 12. Configuration Management

### 12.1 System Configuration

Navigate to **Configuration** to manage:

Category	Settings
<b>General</b>	Platform name, domain, timezone
<b>Email</b>	SMTP settings, templates
<b>Security</b>	Password policy, MFA settings
<b>API</b>	Rate limits, CORS settings
<b>Features</b>	Feature flags

### 12.2 Tenant Overrides

Allow tenant-specific configuration:

1. Navigate to **Configuration** → **Tenant Overrides**

2. Select tenant
3. Override specific settings
4. Save changes

### 12.3 SSM Parameters

System configuration is stored in AWS SSM:

Parameter	Description
/radiant/prod/database/url	Database connection
/radiant/prod/api/rate-limit	API rate limits
/radiant/prod/features/*	Feature flags

## 13. AI Ethics & Standards

**Location:** Admin Dashboard → Ethics

The Ethics module provides transparency into the ethical principles guiding AI behavior and their source standards.

### 13.1 Standards Tab

View all industry AI ethics frameworks that inform RADIANT's ethical principles:

Standard	Organization	Type	Required
<b>NIST AI RMF 1.0</b>	National Institute of Standards and Technology	Government	
<b>ISO/IEC 42001:2023</b>	International Organization for Standardization	ISO	
<b>EU AI Act</b>	European Union	Government	
<b>IEEE 7000-2021</b>	IEEE	Industry	
<b>OECD AI Principles</b>	OECD	Government	
<b>UNESCO AI Ethics</b>	UNESCO	Government	
<b>ISO/IEC 23894:2023</b>	ISO	ISO	
<b>ISO/IEC 38507:2022</b>	ISO	ISO	

Each standard shows:

- **Full Name:** Complete standard title with version
- **Organization:** Issuing body
- **Type:** Government, ISO, Industry, Academic, Religious
- **Required Badge:** Mandatory for compliance
- **Description:** Summary of the standard's purpose
- **Publication Date:** When the standard was issued
- **External Link:** Direct link to official standard

### 13.2 Principles Tab

View ethical principles with their standard sources:

Principle	Category	Source Standards
Love Others	Love	NIST GOVERN 1.2, ISO 42001 Clause 5.2, Matthew 22:39
Golden Rule	Love	NIST MAP 1.1, EU AI Act Article 9, Matthew 7:12

Principle	Category	Source Standards
Speak Truth	Truth	NIST GOVERN 4.1, ISO 42001 Clause 7.4, EU AI Act Article 13, John 8:32
Show Mercy	Mercy	NIST MEASURE 2.6, EU AI Act Article 14, Matthew 5:7
Serve Humbly	Service	NIST GOVERN 1.1, ISO 42001 Clause 5.1, Mark 10:45
Avoid Judgment	Mercy	NIST MAP 2.3, EU AI Act Article 10, Matthew 7:1
Care for Vulnerable	Service	NIST MAP 1.5, EU AI Act Article 7, ISO 42001 Clause 8.4, Matthew 25:40

Each principle displays:

- **Teaching:** The principle text
- **Category Badge:** love, mercy, truth, service, humility, peace, forgiveness
- **Derived from / Aligned with:** Standard badges with section references

### 13.3 Violations Tab

Monitor ethical violations in AI responses:

- **Action:** What was evaluated
- **Score:** Ethical score percentage
- **Concerns:** Specific violations detected
- **Guidance:** Recommended corrections
- **Reasoning:** Explanation of the evaluation

### 13.4 Statistics

Summary cards showing:

- **Total Evaluations:** All ethical checks performed
- **Approved:** Passed evaluations
- **Violations:** Failed evaluations
- **Average Score:** Mean ethical score
- **Today:** Evaluations in the last 24 hours

### 13.5 Standard Alignment Levels

Principles map to standards with different alignment levels:

Level	Meaning
<b>derived</b>	Principle was directly derived from this standard
<b>aligned</b>	Principle aligns with this standard's requirements
<b>supports</b>	Principle supports this standard's goals
<b>extends</b>	Principle extends beyond this standard

See AI Ethics Standards Documentation for full details.

## 14. Provider Rejection Handling

**Location:** Think Tank → Notifications (Bell icon)

When AI providers reject prompts based on their policies (not RADIANT's), the system automatically attempts fallback to alternative models.

## 14.1 How It Works

User Prompt → Model Rejects → Check RADIANT Ethics

Our ethics block → Reject with explanation

Our ethics pass → Try fallback models (up to 3 attempts)

Fallback succeeds → Return response (user notified)

All fail → Reject with full explanation

## 14.2 Rejection Types

Type	Description	Auto-Fallback
<b>content_policy</b>	Provider's content policy violation	Yes
<b>safety_filter</b>	Safety/moderation filter triggered	Yes
<b>provider_ethics</b>	Provider's ethical guidelines differ	Yes
<b>capability_mismatch</b>	Model can't handle request type	Yes
<b>context_length</b>	Prompt too long for model	Yes
<b>moderation</b>	Pre-flight moderation blocked	Yes
<b>rate_limit</b>	Rate limiting	Retry

## 14.3 User Notifications

Users see rejection notifications via the bell icon in Think Tank:

- **Unread count badge** - Number of unread notifications
- **Notification panel** - Slides out showing all rejections
- **Suggested actions** - Rephrase, simplify, contact admin
- **Resolution status** - Whether fallback succeeded

## 14.4 Fallback Model Selection

Models are selected for fallback based on:

1. **Rejection rate** - Models with lowest historical rejection rates preferred
2. **Required capabilities** - Must match original model's capabilities
3. **Provider diversity** - Prefer different providers

## 14.5 Model Rejection Statistics

**Location:** Admin Dashboard → Analytics → Model Stats

View per-model:

- Total requests and rejections
- Rejection rate percentage
- Breakdown by rejection type
- Fallback success rate

## 14.6 Database Tables

Table	Purpose
<b>provider_rejections</b>	Track all rejections with fallback chain
<b>rejection_patterns</b>	Learn patterns for smarter fallback
<b>user_rejection_notifications</b>	User-facing notifications
<b>model_rejection_stats</b>	Per-model statistics



## 14.7 Configuration

Setting	Default	Description
MAX_FALLBACK_ATTEMPTS	3	Maximum fallback tries per request
MIN_MODELS_FOR_TASK	2	Minimum capable models required

See Provider Rejection Handling Documentation for full details.

## 14.8 Rejection Analytics

**Location:** Admin Dashboard → Analytics → Rejections

Comprehensive analytics on AI provider rejections to inform ethics policy updates.

### Summary Cards

- **Total Rejections (30d)** - All rejections in last 30 days
- **Fallback Success Rate** - Percentage resolved by fallback models
- **Rejected to User** - Requests that couldn't be completed
- **Flagged Keywords** - Keywords marked for policy review

### By Provider Tab

Column	Description
Provider	Provider ID (openai, anthropic, google, etc.)
Rejections	Total rejections in period
Models	Number of models affected
Unique Prompts	Distinct prompts rejected
Fallback Rate	Success rate of fallback attempts
Rejected to User	Requests that failed all fallbacks
Types	Rejection types (content_policy, safety_filter, etc.)

**Violation Keywords Tab** Track which keywords trigger rejections:

Column	Description
Keyword	The triggering keyword
Category	violence, security, controlled, etc.
Occurrences	Times keyword appeared in rejected prompts
Provider Breakdown	Rejections per provider for this keyword
Status	Flagged, Pre-filter added, or Monitoring

### Actions:

- **Flag for Review** - Mark keyword for policy consideration
- **Add Pre-Filter** - Block prompts containing keyword before sending to AI

**Flagged Prompts Tab** View full content of rejected prompts for investigation:

- Full prompt text (for policy review only)
- Detected violation keywords
- Model and provider that rejected
- Number of times this prompt pattern was rejected
- Suggested actions: Add Pre-Filter, Dismiss

**Policy Review Tab** Recommendations based on rejection patterns:

- High-frequency rejection patterns
- Suggested pre-filters to add
- Keywords to consider blocking

#### Database Tables

Table	Purpose
rejection_analytics	Daily aggregated stats by model/provider/mode
rejection_keyword_stats	Per-keyword occurrence and rejection counts
rejected_prompt_archive	Full prompt content for flagged reviews

#### Using Analytics to Update Ethics Policy

1. **Monitor** → Watch the Keywords tab for high-occurrence terms
  2. **Flag** → Mark suspicious keywords for review
  3. **Investigate** → View full prompts in Flagged Prompts tab
  4. **Decide** → Add pre-filter, add warning, or dismiss
  5. **Implement** → Update RADIANT ethics policy to pre-filter prompts
- 

## 15. Security & Compliance

### 15.1 Security Dashboard

Navigate to **Security** to monitor:

Security Dashboard

Threat Level: Low

Active Threats: 0  
Failed Logins: 23 (last 24h)  
Suspicious IPs: 2 blocked  
MFA Adoption: 94%

Recent Alerts:

Unusual login location - user@acme.com (2h ago)  
Resolved: Brute force attempt blocked (5h ago)  
Resolved: API key rotated - tenant xyz (1d ago)

### 13.2 Anomaly Detection

Automatic detection of:

- Impossible travel (geographic anomalies)
- Session hijacking attempts
- Brute force attacks
- Unusual API patterns

### 15.3 Compliance Dashboard

Navigate to **Compliance** to access five tabs:

Tab	Purpose
<b>Reports</b>	SOC 2, HIPAA, GDPR, ISO 27001 compliance reports
<b>GDPR</b>	Data subject requests and consent management
<b>HIPAA</b>	PHI protection and access controls
<b>Breaches</b>	Data breach incident tracking
<b>Retention</b>	Data retention policy configuration

### 15.4 GDPR Management

The GDPR tab provides:

**Data Subject Requests** (Articles 15-22): | Request Type | Description | Deadline | |-----|-----|-----|  
-----| | Access | Export all user data | 30 days | | Rectification | Correct personal data | 30 days | | Erasure  
| Right to be forgotten | 30 days | | Restriction | Limit processing | 30 days | | Portability | Machine-readable  
export | 30 days | | Objection | Object to processing | 30 days |

#### Processing Requests:

1. Navigate to **Compliance** → **GDPR**
2. View pending requests with deadlines
3. Click **”Process”** to fulfill request
4. System automatically exports/deletes data
5. Request marked complete with audit trail

### 15.5 HIPAA Configuration

The HIPAA tab provides per-tenant settings:

Setting	Default	Description
<b>HIPAA Mode</b>	Off	Enable all HIPAA features
<b>PHI Detection</b>	On	Auto-detect PHI in requests
<b>PHI Encryption</b>	On	Column-level encryption
<b>Enhanced Logging</b>	On	Detailed PHI access logs
<b>MFA Required</b>	On	Require MFA for all users
<b>Session Timeout</b>	15 min	Auto-logout after inactivity
<b>Access Review</b>	90 days	Periodic access review period
<b>PHI Retention</b>	2190 days	6 years per HIPAA
<b>Audit Retention</b>	2555 days	7 years recommended

### 15.6 Data Breach Management

The Breaches tab tracks security incidents:

#### Incident Types:

- Unauthorized access
- Data theft
- Ransomware
- Accidental disclosure
- System breach
- Insider threat

**Notification Requirements:** | Regulation | Notify | Timeline | |-----|-----|-----| | GDPR | DPA  
+ affected users | 72 hours | | HIPAA | HHS + affected individuals | 60 days |

## 15.7 Data Retention Policies

Default retention periods:

Data Type	Retention	Legal Basis	Action
Session Data	90 days	Legitimate Interest	Delete
Usage Analytics	2 years	Legitimate Interest	Anonymize
Audit Logs	7 years	Legal Obligation	Archive
PHI Data	6 years	Legal Obligation	Archive
Billing Records	7 years	Legal Obligation	Archive
GDPR Requests	3 years	Legal Obligation	Archive
Consent Records	7 years	Legal Obligation	Archive

## 15.8 Generating Compliance Reports

1. Navigate to **Compliance** → **Reports**
2. Select framework (SOC 2, HIPAA, GDPR, ISO 27001)
3. Click **"Generate Report"**
4. View compliance score and findings
5. Export to PDF for auditors

## 15.9 Regulatory Standards Registry

**Location:** Admin Dashboard → Security → Reg Standards

Comprehensive registry of all regulatory standards Radiant must comply with.

### Supported Standards (35 Frameworks)

Category	Standards
<b>Data Privacy</b>	GDPR, CCPA, CPRA, LGPD, PIPEDA, APPI, PDPA
<b>Healthcare</b>	HIPAA, HITECH, HITRUST CSF
<b>Security</b>	SOC 2, SOC 1, ISO 27001, ISO 27017, ISO 27018, ISO 27701, CSA STAR, NIST CSF, NIST 800-53, C
<b>Financial</b>	PCI-DSS, SOX, GLBA
<b>Government</b>	FedRAMP, StateRAMP, ITAR, CMMC
<b>AI Governance</b>	EU AI Act, NIST AI RMF, ISO 42001, IEEE 7000
<b>Accessibility</b>	WCAG 2.1, ADA, Section 508
<b>Industry</b>	FERPA, COPPA

**Overview Tab** The Overview tab provides:

- **Summary Cards:** Total standards, requirements, implementation progress
- **Categories Grid:** Standards organized by domain with mandatory counts
- **Priority Standards:** Mandatory frameworks requiring immediate attention
- **Requirements Status:** Not started, in progress, implemented, verified

**Standards Tab** Browse all regulatory frameworks with:

- **Search:** Find standards by code, name, or description
- **Category Filter:** Filter by Data Privacy, Security, Healthcare, etc.

- **Mandatory Filter:** Show only required standards
- **Standard Details:** Click any standard to view full requirements

**My Compliance Tab** Track your tenant's compliance status:

- **Enable/Disable:** Toggle which standards apply to your organization
- **Compliance Score:** Track progress per standard (0-100%)
- **Status:** Not Assessed → Non-Compliant → Partial → Compliant → Certified
- **Audit Tracking:** Last audit date and next scheduled audit

**Requirements Tracker** For each standard, track individual requirements:

Field	Description
<b>Requirement Code</b>	Unique identifier (e.g., GDPR-32, PCI-7)
<b>Title</b>	Short description of the requirement
<b>Control Type</b>	Technical, Administrative, Physical, Procedural
<b>Status</b>	Not Started → In Progress → Implemented → Verified
<b>Owner</b>	Person responsible for implementation
<b>Evidence</b>	Link to compliance evidence
<b>Due Date</b>	Implementation deadline

### Updating Requirement Status

1. Navigate to **Reg Standards** → **Standards**
2. Click on a standard to open details sheet
3. Find the requirement to update
4. Use the status dropdown to change implementation status
5. Changes are automatically saved

## 15.10 Self-Audit & Regulatory Reporting

**Location:** Admin Dashboard → Security → Self-Audit

Automated compliance self-auditing with timestamped pass/fail results for regulatory reporting.

### Running an Audit

1. Navigate to **Security** → **Self-Audit**
2. Click **"Run Audit"** button
3. Select framework (SOC 2, HIPAA, GDPR, ISO 27001, PCI-DSS, or All)
4. Audit executes automatically (~5-30 seconds depending on scope)
5. Results appear in dashboard with pass/fail breakdown

**Dashboard Overview** The dashboard displays:

- **Overall Pass Rate:** Aggregate compliance score across all frameworks
- **Total Checks:** Number of automated compliance checks (45+)
- **Critical Issues:** Failed checks with critical severity requiring immediate attention
- **Framework Scores:** Individual compliance scores per framework with trend indicators

### Framework Compliance Checks

Framework	Checks	Categories
<b>SOC 2</b>	12	Access Control, Data Protection, Audit Logging, Incident Response, Change Management
<b>HIPAA</b>	8	PHI Protection, Access Control, Audit Trail, Data Retention
<b>GDPR</b>	8	Consent Management, Data Subject Rights, Data Processing, Data Transfer, Breach Response
<b>ISO 27001</b>	9	Security Policy, Organization, Access Control, Cryptography, Operations, Incident Management, C
<b>PCI-DSS</b>	8	Network Security, Data Protection, Encryption, Access Control, Authentication, Monitoring, Testin

**Audit History** View historical audit runs with:

- **Timestamp:** Exact date/time of audit execution
- **Framework:** Which standard was audited
- **Type:** Manual, Scheduled, or Triggered
- **Score:** Overall compliance percentage
- **Pass/Fail/Skip:** Breakdown of check results
- **Triggered By:** Admin who initiated the audit

**Viewing Audit Details** Click any audit run to see:

- **Score Overview:** Large score display with pass/fail/skip counts
- **Critical Failures:** Red panel highlighting checks requiring immediate action
- **By Category:** Progress bars showing compliance per category
- **All Results:** Expandable list of every check with status and remediation steps

**Generating Reports** From an audit run details view:

1. Click **"Export PDF Report"** for auditor-ready documentation
2. Click **"View Evidence"** to see raw query results and execution details

**Scheduling Audits** Configure automated audit schedules:

1. Navigate to **Self-Audit → Checks Registry**
2. Select framework to schedule
3. Choose frequency: Daily, Weekly, Monthly, or Quarterly
4. Set notification preferences for failures
5. Audits run automatically at configured times

**API Integration** Integrate audit results with external systems:

```
# Run audit programmatically
POST /api/admin/self-audit/run
{ "framework": "soc2" }

# Fetch latest results
GET /api/admin/self-audit/runs/{runId}

# Generate report
GET /api/admin/self-audit/runs/{runId}/report
```

## 15.11 Compliance Framework Reference

Comprehensive reference for all supported compliance frameworks with implementation details.

### SOC 2 Type II Trust Services Criteria Implementation:

Control	Radiant Implementation	Evidence
<b>CC1.1 - COSO Integrity</b>	Ethics pipeline with content screening	Ethics audit logs
<b>CC2.1 - Security Policy</b>	Tenant-level security configuration	<code>dynamic_config</code> table
<b>CC3.1 - Risk Assessment</b>	Security anomaly detection	<code>security_anomalies</code> table
<b>CC4.1 - Monitoring</b>	Real-time audit logging	DynamoDB + PostgreSQL audit trails
<b>CC5.1 - Logical Access</b>	Cognito User Pools with MFA	AWS CloudTrail
<b>CC5.2 - Authentication</b>	JWT tokens with session management	Token validation logs
<b>CC6.1 - Encryption</b>	AES-256 at rest, TLS 1.3 in transit	AWS KMS key policies
<b>CC6.6 - Change Management</b>	Dual-approval for production changes	<code>approval_requests</code> table
<b>CC7.1 - System Operations</b>	CloudWatch monitoring, automated alerts	AWS CloudWatch dashboards
<b>CC7.2 - Incident Response</b>	Security anomaly alerting	SNS notifications
<b>CC8.1 - Change Management</b>	Database migration approvals	Migration audit logs
<b>CC9.1 - Business Continuity</b>	Multi-AZ deployment, automated backups	Aurora automated backups

#### Automated Audit Checks:

- MFA enforcement for administrators
- Password policy configuration
- Session timeout settings ( 30 minutes)
- RBAC implementation verification
- Encryption at rest confirmation
- Audit log retention ( 7 years)
- Change management process validation

#### HIPAA / HITECH Administrative Safeguards (45 CFR §164.308):

Requirement	Implementation	Configuration
<b>Security Officer</b>	Designated in tenant config	<code>hipaa_config.security_officer</code>
<b>Workforce Training</b>	User acknowledgment tracking	Onboarding audit logs
<b>Access Management</b>	Role-based access with MFA	Cognito + RLS policies
<b>Contingency Plan</b>	Multi-AZ, automated backups	Aurora + S3 replication
<b>Evaluation</b>	Self-audit system	Quarterly compliance audits

#### Technical Safeguards (45 CFR §164.312):

Requirement	Implementation	Evidence
<b>Access Control</b>	Unique user IDs, automatic logoff	Cognito user management
<b>Audit Controls</b>	PHI access logging	<code>phi_access_logs</code> table
<b>Integrity Controls</b>	Column-level encryption	AWS KMS + application encryption
<b>Transmission Security</b>	TLS 1.3 enforced	API Gateway + ALB policies

#### PHI Detection & Protection:

##### PHI Detection Pipeline

```

Request → PHI Scanner → [PHI Detected?]
                        ↓ Yes
                    Encrypt + Log → Sanitize → Continue
                        ↓ No

```

## Continue Processing

### Retention Requirements:

- PHI Data: 6 years (2,190 days)
- Audit Logs: 7 years (2,555 days)
- BAA Records: 6 years from termination

### GDPR (General Data Protection Regulation) Lawful Bases Supported:

Basis	Use Case	Configuration
<b>Consent (Art. 6(1)(a))</b>	Marketing, analytics	Consent tracking UI
<b>Contract (Art. 6(1)(b))</b>	Service delivery	Terms acceptance
<b>Legal Obligation (Art. 6(1)(c))</b>	Audit retention	Automatic retention
<b>Legitimate Interest (Art. 6(1)(f))</b>	Security, fraud prevention	Documented in DPA

### Data Subject Rights Implementation:

Right	Article	API Endpoint	SLA
<b>Access</b>	Art. 15	POST /gdpr/request/access	30 days
<b>Rectification</b>	Art. 16	PATCH /users/{id}	30 days
<b>Erasure</b>	Art. 17	POST /gdpr/request/erasure	30 days
<b>Restriction</b>	Art. 18	POST /gdpr/request/restriction	30 days
<b>Portability</b>	Art. 20	POST /gdpr/request/portability	30 days
<b>Objection</b>	Art. 21	POST /gdpr/request/objection	30 days

### Cross-Border Transfers:

- Standard Contractual Clauses (SCCs) for non-EU transfers
- AWS regions: **eu-west-1**, **eu-central-1** for EU data residency
- Region restrictions configurable per tenant

### Breach Notification:

- 72-hour notification to supervisory authority
- Affected user notification without undue delay
- Breach tracking in **data\_breaches** table

### ISO 27001:2022 Annex A Controls Implementation:

Control	Title	Radiant Implementation
<b>A.5.1</b>	Policies for information security	<b>dynamic_config</b> security policies
<b>A.5.15</b>	Access control	Cognito + RLS + RBAC
<b>A.5.23</b>	Information security for cloud services	AWS security configurations
<b>A.6.1</b>	Screening	Admin approval workflow
<b>A.8.2</b>	Privileged access rights	Super Admin role restrictions
<b>A.8.3</b>	Information access restriction	Tenant isolation (RLS)
<b>A.8.5</b>	Secure authentication	MFA, JWT, session management
<b>A.8.9</b>	Configuration management	Infrastructure as Code (CDK)
<b>A.8.10</b>	Information deletion	Automated retention + deletion
<b>A.8.12</b>	Data leakage prevention	PHI detection, ethics pipeline



Control	Title	Radiant Implementation
<b>A.8.15</b>	Logging	Comprehensive audit logging
<b>A.8.16</b>	Monitoring activities	CloudWatch + anomaly detection
<b>A.8.24</b>	Use of cryptography	AWS KMS, AES-256, TLS 1.3

#### PCI-DSS v4.0 Requirements Mapping:

Req	Description	Implementation
<b>1</b>	Network security controls	VPC security groups, NACLs
<b>2</b>	Secure configurations	CDK hardened templates
<b>3</b>	Protect stored data	AES-256 encryption, tokenization
<b>4</b>	Protect data in transit	TLS 1.3, certificate pinning
<b>5</b>	Malware protection	AWS WAF, Shield
<b>6</b>	Secure systems	Automated patching, code review
<b>7</b>	Restrict access	RBAC, least privilege
<b>8</b>	Identify users	Unique IDs, MFA, password policies
<b>9</b>	Physical access	AWS data center controls
<b>10</b>	Logging & monitoring	CloudTrail, audit logs
<b>11</b>	Security testing	Automated scanning, pen testing
<b>12</b>	Security policies	Documented in admin guide

**Note:** Full PCI-DSS certification requires additional controls beyond Radiant's scope (physical security, organizational policies). Radiant provides the technical controls for SAQ-A or SAQ-D compliance.

#### FedRAMP / StateRAMP Control Families Addressed:

Family	Controls	Radiant Coverage
<b>AC</b>	Access Control	Full
<b>AU</b>	Audit & Accountability	Full
<b>CM</b>	Configuration Management	Full
<b>IA</b>	Identification & Authentication	Full
<b>SC</b>	System & Communications Protection	Full
<b>SI</b>	System & Information Integrity	Full

#### FedRAMP Boundaries:

- Authorization Boundary: AWS GovCloud available
- Data Flow: Documented in System Security Plan
- Interconnections: API Gateway with mutual TLS

#### EU AI Act Compliance Risk Classification:

Category	AI Systems	Radiant Controls
<b>Unacceptable</b>	Social scoring, subliminal manipulation	Blocked by ethics pipeline
<b>High-Risk</b>	Healthcare, legal, employment	Enhanced logging, human oversight
<b>Limited</b>	Chatbots, emotion recognition	Transparency disclosures
<b>Minimal</b>	Spam filters, games	Standard controls

#### Article 9 - Risk Management:

- Domain detection for high-risk use cases
- Ethics screening with configurable policies
- Human-in-the-loop for sensitive decisions

#### Article 13 - Transparency:

- Model disclosure in API responses
- Training data documentation
- Capability limitations documented

#### Article 14 - Human Oversight:

- Admin approval workflows
- Escalation triggers
- Override capabilities

### 15.12 Compliance Checklist Registry

**Location:** Admin Dashboard → Security → Compliance Checklists

The Compliance Checklist Registry provides versioned, interactive checklists linked to regulatory standards with auto-update support.

#### Key Features

- **Versioned Checklists** - Each regulatory standard has versioned checklists
- **Regulatory Linking** - Checklists are linked to **regulatory\_standards** registry
- **Auto-Update** - Service automatically checks for regulatory version updates
- **Per-Tenant Configuration** - Choose version selection mode per standard
- **Progress Tracking** - Track item completion across teams
- **Audit Runs** - Schedule and run formal checklist audits

#### Version Selection Modes

Mode	Behavior
<b>Auto</b> (default)	Automatically uses latest active checklist version
<b>Specific</b>	Use a specific version, updates when newer available
<b>Pinned</b>	Locked to exact version, no automatic updates

**Dashboard Overview** Navigate to **Security** → **Compliance Checklists** to see:

- All regulatory standards with completion progress
- Per-standard checklist status and version info
- Pending regulatory version updates
- Recent audit run history

**Checklist Items** Each checklist item includes:

Field	Description
<b>Item Code</b>	Unique identifier (e.g., SOC2-PRE-001)
<b>Priority</b>	Critical, High, Medium, Low
<b>Evidence Types</b>	Required evidence (document, screenshot, config, log)
<b>API Endpoint</b>	Optional endpoint for automated evidence collection
<b>Automated Check</b>	Link to <b>system_audit_checks</b> for auto-validation
<b>Estimated Time</b>	Minutes to complete the item

Field	Description
<b>Guidance</b>	Detailed instructions for completing the item

### Item Status Tracking

Status	Description
<b>Not Started</b>	Item has not been addressed
<b>In Progress</b>	Work is underway
<b>Completed</b>	Item is done with evidence attached
<b>Not Applicable</b>	Item doesn't apply to this tenant
<b>Blocked</b>	Item cannot proceed (reason required)

**Pre-Built Checklists** Initial checklists are seeded for:

Standard	Version	Categories
<b>SOC 2 Type II</b>	2024.1	Pre-Audit, Documentation, Evidence, Access Control, Change Mgmt, Risk Mgmt, Monitoring
<b>HIPAA</b>	2024.1	Administrative Safeguards, Technical Safeguards, Physical Safeguards
<b>GDPR</b>	2024.1	Data Subject Rights, Processing Records, Security Measures
<b>ISO 27001:2022</b>	2022.1	Organizational, People, Physical, Technological (93 controls)
<b>PCI-DSS</b>	4.0	12 requirement categories

**Auto-Update Service** The checklist registry can automatically check for regulatory updates:

1. **Update Sources** - Configure RSS feeds, APIs, or webhooks per standard
2. **Check Frequency** - Default 24 hours, configurable per source
3. **Update Detection** - Records detected version changes
4. **Processing** - Admin reviews and approves new checklist versions
5. **Notification** - Tenants notified when their effective version changes

**Audit Runs** Start formal checklist reviews:

Run Type	Purpose
<b>Manual</b>	Ad-hoc review
<b>Scheduled</b>	Periodic compliance check
<b>Pre-Audit</b>	Preparation before external audit
<b>Certification</b>	Formal certification attempt

**API Endpoints** Base Path: `/api/admin/compliance/checklists`

### Dashboard & Configuration:

Endpoint	Method	Description
<code>/dashboard</code>	GET	Full dashboard with progress per standard
<code>/config</code>	GET	All tenant checklist configurations
<code>/config/:standardId</code>	GET/PUT	Get/set version selection for a standard
<code>/config/:standardId/effective-version</code>	GET	Get effective version for tenant

**Versions & Items:**

Endpoint	Method	Description
/versions	GET/POST	List versions for standard / Create version
/versions/latest	GET	Get latest version for standard code
/versions/:id	GET	Get specific version details
/versions/:id/set-latest	POST	Set version as latest
/versions/:id/categories	GET/POST	List/create categories
/versions/:id/items	GET/POST	List/create checklist items

### Progress & Audit:

Endpoint	Method	Description
/progress/:versionId	GET	Get tenant progress for version
/progress/items/:itemId	PUT	Update item progress/status
/audit-runs	GET/POST	List history / Start new audit run
/audit-runs/:id/complete	PUT	Complete an audit run

### Auto-Updates:

Endpoint	Method	Description
/updates/pending	GET	Get pending regulatory updates
/updates	POST	Record a version update
/updates/:id/process	PUT	Process (approve/reject) an update
/updates/check/:standardId	POST	Check sources for updates

### Database Tables

Table	Purpose
compliance_checklist_versions	Versioned checklists per standard
compliance_checklist_categories	Categories within a checklist
compliance_checklist_items	Individual checklist items
tenant_checklist_config	Per-tenant version selection
tenant_checklist_progress	Item completion tracking
checklist_audit_runs	Audit run history
regulatory_version_updates	Detected regulatory updates
checklist_update_sources	Auto-update source configuration

### Quick Reference Checklist For quick pre-audit preparation:

- ☐ Run self-audit for all frameworks (/compliance/self-audit)
- ☐ Export audit logs for review period
- ☐ Generate compliance reports (PDF)
- ☐ Review critical findings and remediation status
- ☐ Verify all evidence artifacts are accessible
- ☐ Confirm data retention policies are enforced

### Documentation Required

Document	Location	Purpose
System Security Plan	/docs/SYSTEM-SECURITY-PLAN.md	Architecture overview
Data Flow Diagram	Admin Dashboard → Compliance	Data processing flows
Access Control Matrix	/docs/ACCESS-CONTROL-MATRIX.md	Role permissions
Incident Response Plan	/docs/INCIDENT-RESPONSE.md	Breach procedures
Business Continuity Plan	/docs/BUSINESS-CONTINUITY.md	Disaster recovery

### API Request/Response Examples Get Dashboard Data:

GET /api/admin/compliance/checklists/dashboard

Headers: x-tenant-id: your-tenant-id

Response:

```
{
  "standards": [
    {
      "id": "standard-123",
      "code": "SOC2",
      "name": "SOC 2 Type II",
      "latestVersion": "2024.1",
      "effectiveVersionId": "version-abc",
      "completionPercentage": 66.7,
      "itemsCompleted": 12,
      "totalItems": 18
    }
  ],
  "pendingUpdates": 0,
  "recentAuditRuns": []
}
```

### Get Versions for Standard:

GET /api/admin/compliance/checklists/versions?standardId=standard-123

Response:

```
{
  "versions": [
    {
      "id": "version-abc",
      "version": "2024.1",
      "title": "SOC 2 Type II Pre-Audit Checklist",
      "versionDate": "2024-01-01",
      "isLatest": true,
      "isActive": true,
      "categoriesCount": 7,
      "itemsCount": 18
    }
  ]
}
```

### Create Checklist Version:

POST /api/admin/compliance/checklists/versions

Content-Type: application/json

```
{
  "standardId": "standard-123",
  "version": "2025.1",
  "title": "SOC 2 Type II Pre-Audit Checklist 2025",
  "description": "Updated for 2025 AICPA guidance",
  "versionDate": "2025-01-01"
}
```

Response: 201 Created

```
{
  "id": "new-version-id",
  "version": "2025.1",
  "title": "SOC 2 Type II Pre-Audit Checklist 2025",
  "isLatest": false,
  "isActive": true
}
```

### Get Items with Progress:

GET /api/admin/compliance/checklists/versions/{versionId}/items  
 Headers: x-tenant-id: your-tenant-id

Response:

```
{
  "items": [
    {
      "id": "item-001",
      "itemCode": "SOC2-PRE-001",
      "title": "Confirm audit dates",
      "description": "Schedule dates with external auditor",
      "categoryCode": "pre_audit",
      "priority": "critical",
      "isRequired": true,
      "estimatedMinutes": 15,
      "evidenceTypes": ["document", "attestation"],
      "status": "completed",
      "completedAt": "2024-01-15T10:30:00Z",
      "completedBy": "user-123"
    }
  ]
}
```

### Update Item Progress:

PUT /api/admin/compliance/checklists/progress/items/{itemId}  
 Content-Type: application/json  
 Headers: x-tenant-id: your-tenant-id

```
{
  "status": "completed",
  "notes": "Verified with auditor on call",
  "evidenceUrls": [
    "https://storage.example.com/evidence/audit-confirmation.pdf"
  ]
}
```

Response:

```
{ "success": true }
```

### Set Tenant Version Configuration:

PUT /api/admin/compliance/checklists/config/{standardId}

Content-Type: application/json

Headers: x-tenant-id: your-tenant-id

```
{
  "versionSelection": "specific",
  "selectedVersionId": "version-abc",
  "autoUpdateEnabled": false,
  "notificationOnUpdate": true
}
```

Response:

```
{
  "tenantId": "your-tenant-id",
  "standardId": "standard-123",
  "versionSelection": "specific",
  "selectedVersionId": "version-abc",
  "autoUpdateEnabled": false,
  "notificationOnUpdate": true
}
```

### Start Audit Run:

POST /api/admin/compliance/checklists/audit-runs

Content-Type: application/json

Headers: x-tenant-id: your-tenant-id

```
{
  "versionId": "version-abc",
  "runType": "pre_audit",
  "notes": "Preparing for Q1 external audit"
}
```

Response: 201 Created

```
{
  "id": "run-123",
  "tenantId": "your-tenant-id",
  "versionId": "version-abc",
  "runType": "pre_audit",
  "status": "in_progress",
  "startedAt": "2024-01-15T10:00:00Z",
  "triggeredBy": "user-123"
}
```

### Complete Audit Run:

PUT /api/admin/compliance/checklists/audit-runs/{runId}/complete

Content-Type: application/json

```
{
  "status": "completed",
  "score": 95,
}
```

```

    "findings": [
      "Minor documentation gap in CC5.2 - recommend update before external audit"
    ],
    "recommendations": [
      "Schedule follow-up review for access control documentation"
    ]
  }
}

```

Response:

```

{
  "id": "run-123",
  "status": "completed",
  "completedAt": "2024-01-16T14:30:00Z",
  "score": 95,
  "findings": ["..."]
}

```

## Query Parameters Reference

Endpoint	Parameter	Type	Description
/versions	standardId	UUID	Required. Filter by standard
/versions/latest	standardCode	String	Required. Standard code (e.g., SOC2)
/versions/:id/items	categoryCode	String	Optional. Filter by category
/audit-runs	limit	Integer	Max results (default: 20)
/config	-	-	Returns all configs for tenant

## Error Responses

Status	Error	Description
400	Bad Request	Missing required field or invalid format
401	Unauthorized	Missing or invalid authentication
403	Forbidden	Tenant ID mismatch or insufficient permissions
404	Not Found	Resource doesn't exist
500	Server Error	Internal error (check logs)

Example error response:

```

{
  "error": "standardId, version, and title required"
}

```

## Evidence Collection API

*# Export audit logs for date range*

GET /api/admin/audit-logs/export?start=2024-01-01&end=2024-12-31

*# Export compliance report*

GET /api/admin/self-audit/runs/{runId}/report

*# Export user access logs*

GET /api/admin/users/access-logs/export



```
# Get checklist progress
GET /api/admin/compliance/checklists/progress/{versionId}

# Start pre-audit checklist run
POST /api/admin/compliance/checklists/audit-runs
Body: { "versionId": "...", "runType": "pre_audit" }
```

### 15.13 Administrator Emergency Protocols

Critical procedures for security incidents and platform emergencies.

**Emergency Tenant Suspension Scenario:** GuardDuty detects malicious activity originating from a specific tenant (e.g., launching an outbound DDoS attack, data exfiltration attempt).

#### Procedure:

1. **Do NOT** shut down the entire system
2. Access **Admin Dashboard** → **Tenants**
3. Locate the tenant by ID or name
4. Click **"Suspend Tenant"** (requires MFA confirmation)
5. Enter suspension reason for audit log

#### System Response:

- API Gateway Authorizer immediately rejects all JWTs containing that `tenant_id`
- Tenant is effectively quarantined from compute layer immediately
- All active sessions are invalidated
- Scheduled jobs for that tenant are paused
- Alert sent to security team via SNS

#### Post-Incident:

1. Investigate root cause using X-Ray traces filtered by `tenant_id`
2. Collect forensic evidence from CloudWatch Logs
3. Determine if data breach occurred (breach notification may be required)
4. Decide: remediate and restore, or permanent termination
5. Document incident in compliance system

**Key Rotation Emergency Scenario:** Compromise of an administrative credential or API key.

#### Procedure:

Compromise Type	Action
IAM User Keys	Rotate immediately via AWS Console
Database Credentials	Trigger Aurora Master Password Rotation via Secrets Manager
Tenant API Keys	Revoke all API keys via Admin Dashboard → Tenant → API Keys
KMS Key Suspected	Schedule key rotation (cannot immediately delete due to data access)

#### System Response:

- Secrets Manager automatically updates the secret
- Database connections automatically restart with new credentials
- No code deployment required
- Audit trail recorded in CloudTrail

**Mass Incident Response Scenario:** Platform-wide security incident (e.g., zero-day vulnerability exploitation).

**Procedure:**

1. **Activate Incident Response Team** - PagerDuty escalation
2. **Enable WAF Emergency Rules** - Block suspicious patterns
3. **Freeze Deployments** - Halt all CI/CD pipelines
4. **Capture Evidence** - Export CloudTrail, CloudWatch, GuardDuty findings
5. **Patch/Mitigate** - Apply hotfix or WAF rule
6. **Gradual Restoration** - Enable services tenant by tenant
7. **Post-Mortem** - Root cause analysis within 72 hours
8. **Communication** - Status page updates, customer notification if data affected

**Data Breach Response Timeline**

Regulation	Notification Deadline	Recipient
GDPR	72 hours	Supervisory Authority
HIPAA	60 days	HHS, affected individuals
State Laws	Varies (often 72 hours)	State AG, affected individuals

**Breach Response Steps:**

1. **Contain** - Suspend affected tenant(s), revoke compromised credentials
2. **Assess** - Determine scope: which tenants, which data, how many records
3. **Preserve** - Forensic copy of affected systems
4. **Report** - Legal team notified within 4 hours
5. **Notify** - Regulatory notifications per timeline
6. **Remediate** - Fix vulnerability, rotate credentials
7. **Document** - Complete incident report for compliance

**Recovery Objectives**

Data Class	RPO (Recovery Point)	RTO (Recovery Time)	Backup Strategy
<b>Critical</b> (User accounts, tenant config)	1 hour	1 hour	Continuous replication + hourly snapshots
<b>High</b> (Conversations, messages)	4 hours	4 hours	4-hour snapshots
<b>Medium</b> (Preferences, cached data)	24 hours	8 hours	Daily snapshots
<b>Low</b> (Analytics, aggregates)	7 days	24 hours	Weekly snapshots
<b>Audit Logs</b>	0 (immutable)	4 hours	S3 Object Lock + cross-region

### 15.14 Security Monitoring Schedules

**Location:** Admin Dashboard → Security → Schedules

Full-featured EventBridge schedule management with templates, notifications, and webhooks.

**Available Schedules**

Schedule	Default	Purpose
<b>Drift Detection</b>	Daily 00:00 UTC	Monitors model output distribution changes
<b>Anomaly Detection</b>	Hourly	Behavioral scans for suspicious patterns
<b>Classification Review</b>	Every 6 hours	Aggregates classification statistics

Schedule	Default	Purpose
<b>Weekly Security Scan</b>	Sunday 02:00 UTC	Comprehensive security audit
<b>Weekly Benchmark</b>	Saturday 03:00 UTC	TruthfulQA and factual accuracy tests

## Managing Schedules

1. Navigate to **Security** → **Schedules**
2. View all schedules with current configuration and next execution time
3. Toggle **Enable/Disable** to control schedule
4. Click **Configure** to modify cron expression (with real-time preview)
5. Click **Run Now** to trigger immediate execution
6. Click **Test Run** for dry-run validation without affecting data

## Bulk Operations

- **Enable All:** Enable all schedules at once
- **Disable All:** Disable all schedules (e.g., for maintenance)

**Schedule Templates** Apply pre-configured schedule sets:

Template	Description
<b>Production (Conservative)</b>	Fewer checks, less resource usage
<b>Development (Aggressive)</b>	Frequent checks for dev environments
<b>Minimal</b>	Weekly only, for low-traffic tenants

To apply a template:

1. Go to **Templates** tab
2. Review template settings
3. Click **Apply** to update all schedules

**Notifications** Configure alerts for schedule executions:

1. Go to **Notifications** tab
2. Enable notifications
3. Configure:
  - **SNS Topic ARN:** For AWS SNS notifications
  - **Slack Webhook URL:** For Slack channel alerts
  - **Notify on Success:** Alert when schedules complete
  - **Notify on Failure:** Alert when schedules fail (recommended)

**Webhooks** Send execution results to external services:

1. Go to **Webhooks** tab
2. Enter webhook URL
3. Click **Add Webhook**

Webhook events:

- **execution.completed** - Schedule finished successfully
- **execution.failed** - Schedule encountered errors

Webhook payload:

```
{
  "event": "execution.completed",
  "payload": {
    "scheduleType": "drift_detection",
    "status": "completed",
    "itemsProcessed": 150,
    "itemsFlagged": 3,
    "executionTimeMs": 4523
  },
  "timestamp": "2024-12-29T12:00:00Z"
}
```

**Cron Expression Format** Uses AWS EventBridge cron format: Minutes Hours Day-of-month Month Day-of-week Year

**Common Presets:** | Preset | Cron Expression | Description | |-----|-----|-----| | Hourly | 0 \* \* \* ? \* | Every hour | | Every 6 hours | 0 0,6,12,18 \* \* ? \* | 4 times daily | | Daily midnight | 0 0 \* \* ? \* | Once per day | | Weekly Sunday 2AM | 0 2 ? \* SUN \* | Weekly on Sunday |

The UI shows human-readable descriptions and next 5 execution times for any cron expression.

**Execution History** View past execution results including:

- **Status:** Running, Completed, Failed
- **Duration:** Execution time in seconds
- **Items Processed:** Number of items checked
- **Items Flagged:** Security issues found
- **Errors:** Any errors encountered

**Audit Log** All schedule changes are logged with:

- Timestamp and user who made change
- Previous and new cron expression
- Enable/disable actions
- Reason for change (optional)

**API Endpoints Core:** | Endpoint | Method | Description | |-----|-----|-----| | /api/admin/security/schedules | GET | Get all schedule config | | /api/admin/security/schedules/dashboard | GET | Dashboard with stats | | /api/admin/security/schedules/{type} | PUT | Update schedule | | /api/admin/security/schedules/{type}/enable | POST | Enable schedule | | /api/admin/security/schedules/{type}/disable | POST | Disable schedule | | /api/admin/security/schedules/{type}/run-now | POST | Trigger manual run (with optional dryRun) |

**Templates:** | Endpoint | Method | Description | |-----|-----|-----| | /api/admin/security/schedules/templates | GET | List templates | | /api/admin/security/schedules/templates | POST | Create template | | /api/admin/security/schedules/templates/{id}/apply | POST | Apply template | | /api/admin/security/schedules/templates/{id} | DELETE | Delete template |

**Notifications & Webhooks:** | Endpoint | Method | Description | |-----|-----|-----| | /api/admin/security/schedules/notifications | GET/PUT | Notification config | | /api/admin/security/schedules/webhooks/{id} | DELETE | Delete webhook |

**Utilities:** | Endpoint | Method | Description | |-----|-----|-----| | /api/admin/security/schedules/parse-cron | POST | Parse cron with preview | | /api/admin/security/schedules/bulk/enable | POST | Enable all schedules | | /api/admin/security/schedules/bulk/disable | POST | Disable all schedules |

/api/admin/security/schedules/presets | GET | Get cron presets | | /api/admin/security/schedules/executions  
| GET | Get execution history | | /api/admin/security/schedules/audit | GET | Get audit log |

---

## 14. Cost Analytics

### 14.1 Cost Dashboard

Navigate to **Cost** to view:

Cost Analytics		Period: Last 30 Days
Total Spend:	\$12,456.78	(+12% vs last month)
Projected:	\$14,200.00	(this month)
By Provider:		
OpenAI:	\$6,234.56	(50%)
Anthropic:	\$3,456.78	(28%)
Self-hosted:	\$1,234.56	(10%)
Other:	\$1,530.88	(12%)
AI Recommendations:		
Switch 23% of GPT-4 calls to GPT-4-mini (save \$890/mo)		
Enable caching for repeated queries (save \$340/mo)		

### 14.2 Cost Alerts

Configure alerts:

- Daily budget exceeded
- Weekly spend spike
- Per-tenant limits
- Per-model thresholds

### 14.3 Cost Optimization

Review AI-powered recommendations:

1. Navigate to **Cost** → **Insights**
  2. Review suggestions
  3. Click **"Apply"** to implement (requires approval)
  4. Track savings over time
- 

## 15. Revenue Analytics

### 15.1 Revenue Dashboard

Navigate to **Revenue** to view gross revenue, COGS, and profit:

Revenue Analytics	Period: Last 30 Days
-------------------	----------------------

Gross Revenue:	\$29,450.00	(+12.5% vs last period)
Total COGS:	\$14,150.00	
Gross Profit:	\$15,300.00	(+15.2% vs last period)
Gross Margin:	51.95%	

Revenue Breakdown:

Subscriptions:	\$15,000.00 (50.9%)
Credit Purchases:	\$2,500.00 (8.5%)
AI Markup (External):	\$8,750.00 (29.7%)
AI Markup (Self-Hosted):	\$3,200.00 (10.9%)

Note: Marketing, sales, G&A costs not included (COGS only)

15.2 Cost Breakdown (COGS)

View infrastructure and provider costs:

Category	Description	Example Services
AWS Compute	Compute infrastructure	EC2, SageMaker, Lambda
AWS Storage	Storage services	S3, EBS
AWS Network	Data transfer	API Gateway, CloudFront
AWS Database	Database services	Aurora, DynamoDB
External AI	Provider costs	OpenAI, Anthropic APIs
Platform Fees	Payment processing	Stripe fees

15.3 Revenue by Model

View per-model profitability:

Model	Provider Cost	Customer Charge	Markup	Requests
gpt-4o	\$500.00	\$650.00	30%	12,345
claude-3.5-sonnet	\$300.00	\$390.00	30%	8,901
Self-hosted Llama	\$100.00	\$175.00	75%	45,678

15.4 Accounting Exports

Export revenue data for accounting software:

1. Click **Export** dropdown
2. Select format:
  - **CSV**: Summary for spreadsheets
  - **JSON**: Full details for integrations
  - **QuickBooks IIF**: Direct import to QuickBooks
  - **Xero CSV**: Import to Xero
  - **Sage CSV**: Import to Sage
3. Configure date range
4. Download file

QuickBooks Integration:

- Import via File → Utilities → Import → IIF Files

- Creates General Journal entries
- Requires matching account names

See Revenue Analytics Documentation for full details.

## 16. SaaS Metrics Dashboard

### 16.1 Overview

Navigate to **SaaS Metrics** for a comprehensive view of business health:

SaaS Metrics Dashboard		Period: Last 30 Days
MRR: \$89,500 +12.5%	ARR: \$1,074,000 +15.2%	Gross Margin: 53.1% 53%
Customers: 342 +5.8%	Churn: 2.3% Target <2%	LTV:CAC: 6.98x Healthy

### 16.2 Key Metrics

Metric	Description	Healthy Range
<b>MRR</b>	Monthly Recurring Revenue	Growing month-over-month
<b>ARR</b>	Annual Recurring Revenue	$\text{MRR} \times 12$
<b>Gross Margin</b>	$(\text{Revenue} - \text{COGS}) / \text{Revenue}$	> 50%
<b>Churn Rate</b>	Customers lost / Total	< 3%
<b>LTV:CAC</b>	Lifetime Value / Acquisition Cost	> 3:1

### 16.3 Dashboard Tabs

1. **Overview:** Revenue trends, top tenants, cost breakdown
2. **Revenue:** MRR movement, revenue by product/tier
3. **Costs:** Cost distribution, COGS breakdown
4. **Customers:** Growth trends, churn analysis
5. **Models:** Per-model profitability analysis

### 16.4 Exporting Reports

Export data for Excel or accounting:

1. Click **Export** dropdown
2. Select format:
  - **Excel (CSV):** Full metrics in spreadsheet format
  - **JSON:** Structured data for integrations
3. File downloads with period and date in filename

**Export includes:**

- Revenue summary (MRR, ARR, Gross Profit)
- Cost breakdown by category
- Customer metrics (total, new, churned)

- Unit economics (ARPU, LTV, CAC)
- Top tenants with details
- Model performance metrics
- Daily trend data

See SaaS Metrics Dashboard Documentation for full details.

## 17. A/B Testing & Experiments

### 16.1 Experiment Dashboard

Navigate to **Experiments** to manage:

Experiment	Status	Variants	Sample Size
Model routing v2	Running	3	45,234
Prompt optimization	Running	2	12,456
Temperature test	Completed	4	89,123

### 15.2 Creating an Experiment

1. Click "+ New Experiment"
2. Configure:
  - **Name:** Descriptive name
  - **Hypothesis:** What you're testing
  - **Variants:** Control + treatments
  - **Traffic Split:** Percentage per variant
  - **Success Metric:** What to measure
3. Set targeting rules
4. Start experiment

### 15.3 Statistical Analysis

View results with:

- Conversion rates per variant
- Statistical significance (p-value)
- Confidence intervals
- Sample size recommendations

## 16. Audit & Monitoring

### 16.1 Audit Logs

Navigate to **Audit** to view all actions:

Column	Description
<b>Timestamp</b>	When action occurred
<b>Actor</b>	Who performed action
<b>Action</b>	What was done
<b>Resource</b>	What was affected
<b>IP Address</b>	Source IP
<b>Details</b>	Additional context



Column	Description
--------	-------------

## 16.2 Log Filtering

Filter by:

- Date range
- Actor (user/admin)
- Action type
- Resource type
- Severity level

## 16.3 Log Export

Export logs for compliance:

1. Set filter criteria
2. Click **"Export"**
3. Choose format (CSV/JSON)
4. Download file

## 16.4 Real-Time Monitoring

Navigate to **Monitoring** for:

- Live request stream
- Error rate graphs
- Latency percentiles
- Active users count

---

## 17. Database Migrations

### 17.1 Migration Workflow

RADIANT uses dual-admin approval for production migrations:

1. **Submit:** Admin submits migration
2. **Review:** Second admin reviews
3. **Approve:** Second admin approves
4. **Execute:** Migration runs
5. **Verify:** Automatic verification

### 17.2 Pending Migrations

Navigate to **Migrations** to see:

Database Migrations

Pending Approval:

```
#045 - Add user preferences table
Submitted by: alice@company.com (2 hours ago)
[View SQL] [Approve] [Reject]
```

Recent Migrations:

- #044 - Cost tracking tables (applied 2024-12-24)
- #043 - Experiment framework (applied 2024-12-20)
- #042 - Security anomalies (applied 2024-12-15)

### 17.3 Approving Migrations

1. Review the SQL in **"View SQL"**
  2. Check for potential issues
  3. Click **"Approve"** or **"Reject"**
  4. Add comment explaining decision
- 

## 18. API Management

### 18.1 API Keys

Manage platform API keys:

1. Navigate to **Settings** → **API Keys**
2. View existing keys
3. Create new keys with scopes
4. Revoke compromised keys

### 18.2 Rate Limiting

Configure rate limits:

Level	Default	Configurable
<b>Global</b>	10,000/min	Yes
<b>Per-Tenant</b>	1,000/min	Yes
<b>Per-User</b>	100/min	Yes
<b>Per-Key</b>	60/min	Yes

### 18.3 Webhooks

Configure outgoing webhooks:

1. Navigate to **Settings** → **Webhooks**
  2. Add webhook URL
  3. Select events to send
  4. Test webhook
  5. Enable webhook
- 

## 19. Anti-Patterns and Prohibited Practices

**CRITICAL:** These patterns represent security vulnerabilities or architectural anti-patterns that must NEVER be used in RADIANT.

## 19.1 Critical Security Anti-Patterns

NEVER DO	INSTEAD DO	Risk
Trust <code>tenant_id</code> from request body	Extract from JWT claims only	Identity spoofing
Store <code>tenant_id</code> in URLs	Use headers/JWT claims	Information disclosure in logs
Use <code>dynamodb:*</code> wildcard permissions	Specify exact actions + Leading Keys	Privilege escalation
Deploy Lambda code without signing	Use AWS Signer verification	Supply chain attack
Allow manual AWS Console changes in prod	Enforce CDK/Terraform only via SCPs	Configuration drift
Use shared <code>/tmp</code> for cross-invocation data	Use DynamoDB/S3 for persistence	Data leakage
Log full request/response bodies	Sanitize sensitive data	PHI/PII exposure
Store credentials in code/environment variables	Use Secrets Manager	Credential theft
Use single KMS key for all tenants	Per-tenant CMK (Tier 3+)	Blast radius on key compromise
Delete audit logs	S3 Object Lock Compliance Mode	Repudiation, compliance failure
Return detailed errors to clients	Log internally, generic response to client	Information disclosure

## 19.2 Architectural Anti-Patterns

AVOID	BETTER APPROACH
Monolithic TMS + Admin Dashboard	Separate Control Plane from Application
Application-layer-only isolation	Database-engine-enforced isolation (RLS, Leading Keys)
Static asset inventory spreadsheets	AWS Config continuous discovery
Point-in-time security scans	Continuous monitoring (GuardDuty + Inspector)
Generic error messages everywhere	Detailed internal logs, generic external responses
Same IAM role for all functions	Function-specific least-privilege roles
Trusting user input for SQL queries	Parameterized queries + RLS

## 19.3 Code Anti-Patterns

```
// WRONG: Tenant ID from untrusted source
const tenantId = event.body.tenant_id;
const tenantId = event.queryStringParameters.tenant_id;
const tenantId = event.headers['X-Tenant-ID'];

// CORRECT: Tenant ID from verified JWT
const tenantId = event.requestContext.authorizer?.claims?.['custom:tenant_id'];

// WRONG: Building queries without RLS context
const result = await db.query(
  `SELECT * FROM conversations WHERE tenant_id = $1`,
  [tenantId]
);

// CORRECT: Set context, let RLS handle isolation
await db.query(`SET app.current_tenant_id = $1`, [tenantId]);
const result = await db.query(`SELECT * FROM conversations`);
// RLS policy automatically filters by tenant_id

// WRONG: Detailed errors to client
return {
  statusCode: 500,
  body: JSON.stringify({
    error: error.message,
  })
}
```

```

        stack: error.stack,
        query: sql,
    }},
};

// CORRECT: Generic error, detailed internal log
logger.error('Query failed', { error, sql, tenantId, traceId });
return {
    statusCode: 500,
    body: JSON.stringify({
        error: 'Internal Server Error',
        requestId: context.awsRequestId
    }),
};

```

## 19.4 STRIDE Threat Model Reference

Threat	Serverless Attack Vector	RADIANT Mitigation
<b>S - Spoofing</b>	Forged JWT, misconfigured authorizer, tenant_id in request body	Cognito + API Gateway
<b>T - Tampering</b>	S3 code bucket compromise, dependency injection, supply chain attack	AWS Signer code signing
<b>R - Repudiation</b>	Lost ephemeral logs, deleted audit trail	CloudWatch → S3
<b>I - Information Disclosure</b>	Shared Lambda /tmp, memory leakage, verbose errors	Lambda Tenant Isolation
<b>D - Denial of Service</b>	Wallet DoS, noisy neighbor, partition exhaustion	Per-tenant API Gateway
<b>E - Elevation of Privilege</b>	Over-permissive IAM roles, role assumption exploits	Least privilege + IAM

## 19.5 Pre-Deployment Security Checklist

- ☐ Lambda Tenant Isolation Mode enabled for all functions
- ☐ RLS policies applied to all tenant-scoped tables
- ☐ DynamoDB Leading Keys IAM conditions configured
- ☐ AWS Signer code signing enabled and enforced
- ☐ S3 Object Lock enabled on audit bucket (Compliance Mode)
- ☐ GuardDuty Lambda Protection enabled
- ☐ Inspector Lambda scanning enabled
- ☐ X-Ray tracing enabled for all functions
- ☐ Per-tenant API Gateway usage plans configured
- ☐ WAF rules deployed and active
- ☐ KMS CMK created for Tier 3+ tenants
- ☐ Permission boundaries applied to all Lambda roles
- ☐ Secrets Manager used for all credentials (no env vars)
- ☐ CloudTrail enabled with log file validation

## 20. Troubleshooting

### 20.1 Common Issues

#### High Error Rate

1. Check **Providers** for unhealthy providers
2. Review **Audit** logs for patterns
3. Check **Monitoring** for load spikes
4. Verify API key validity

## Slow Response Times

1. Check provider latency in **Providers**
2. Review model selection in **Orchestration**
3. Check for cold-start issues (self-hosted)
4. Verify database performance

## Authentication Failures

1. Check user status in **Users**
2. Verify MFA configuration
3. Review **Audit** logs for login attempts
4. Check for IP blocks in **Security**

## 20.2 Support Resources

Resource	Description
<b>Documentation</b>	This guide + online docs
<b>Status Page</b>	status.radiant.example.com
<b>Support Email</b>	support@radiant.example.com
<b>Emergency</b>	+1-555-RADIANT

## 20.3 Log Locations

Service	Log Group
API Gateway	/aws/apigateway/radiant
Lambda	/aws/lambda/radiant-*
Admin Dashboard	/aws/cloudfront/admin
Database	/aws/rds/cluster/radiant

# 20. Delight System Administration

The Delight System provides personality, achievements, and engagement features for Think Tank users.

## 20.1 Accessing Delight Admin

Navigate to **Think Tank** → **Delight** in the admin sidebar.

## 20.2 Dashboard Overview

The Delight dashboard shows:

Metric	Description
<b>Messages Shown</b>	Total delight messages displayed
<b>Achievements Unlocked</b>	Total achievements earned by users
<b>Easter Eggs Found</b>	Hidden features discovered
<b>Active Users</b>	Users with Delight enabled

## 20.3 Managing Categories

Toggle entire categories on/off:

Category	Purpose
Domain Loading	Messages while loading domain expertise
Domain Transition	Messages when switching topics
Time Awareness	Time-of-day contextual messages
Model Dynamics	Messages about AI consensus/disagreement
Complexity Signals	Feedback on query complexity
Synthesis Quality	Post-response quality indicators
Achievements	Milestone celebrations
Wellbeing	Break/health reminders
Easter Eggs	Hidden features
Sounds	Audio feedback

## 20.4 Managing Messages

- **Create:** Add new delight messages with targeting options
- **Edit:** Modify text, triggers, and display settings
- **Delete:** Remove messages (soft delete)
- **Toggle:** Enable/disable individual messages

### Message Targeting Options

Option	Values
Injection Point	pre_execution, during_execution, post_execution
Trigger Type	domain_loading, time_aware, model_dynamics, etc.
Domain Families	science, humanities, creative, technical, etc.
Time Contexts	morning, afternoon, evening, night, weekend
Display Style	subtle, moderate, expressive

## 20.5 Statistics Dashboard

Access detailed usage statistics at **Delight → View Statistics:**

- **Weekly Trends:** 12-week activity history
- **Top Messages:** Most-shown messages with engagement data
- **Achievement Stats:** Unlock rates, time-to-unlock averages
- **Easter Egg Stats:** Discovery rates by egg
- **User Engagement:** Leaderboard by achievement points

## 20.6 Managing Achievements

Configure achievement unlock criteria:

Setting	Description
Threshold	Number required to unlock
Rarity	common, uncommon, rare, epic, legendary
Points	Score value for leaderboards
Hidden	Only visible after unlock

## 20.7 Managing Easter Eggs

Configure hidden features:

Setting	Description
Trigger Type	key_sequence, text_input, time_based, random
Trigger Value	The activation pattern
Effect Type	mode_change, visual_transform, sound_play
Duration	How long the effect lasts

## 20.8 API Endpoints

Endpoint	Method	Description
/api/admin/delight/dashboard	GET	Dashboard data
/api/admin/delight/statistics	GET	Detailed statistics
/api/admin/delight/categories/:id	PATCH	Toggle category
/api/admin/delight/messages	POST	Create message
/api/admin/delight/messages/:id	PUT/DELETE	Update/delete
/api/admin/delight/user-engagement	GET	User leaderboard

## Appendix: Quick Reference

### Keyboard Shortcuts

Shortcut	Action
G + D	Go to Dashboard
G + T	Go to Tenants
G + M	Go to Models
G + B	Go to Billing
G + A	Go to Audit
?	Show shortcuts

### Status Indicators

Icon	Meaning
	Healthy/Success
	Warning
	Error/Failed
	In Progress
	Disabled/Pending

## 16. Adaptive Storage Configuration

**Location:** Admin Dashboard → Settings → Storage

Configure storage backends per deployment tier to optimize costs.

## 16.1 Storage Types

Type	Use Case	Monthly Cost
<b>Fargate PostgreSQL</b>	Tier 1-2 (dev/startup)	\$5-50
<b>Aurora Serverless v2</b>	Tier 3-5 (production)	\$100-2500
<b>DynamoDB</b>	Simple key-value workloads	Pay per request

## 16.2 Default Configuration

Tier	Default Storage	Reason
1 (SEED)	Fargate PostgreSQL	Low cost for dev
2 (STARTUP)	Fargate PostgreSQL	Cost-effective for small prod
3 (GROWTH)	Aurora Serverless	Auto-scaling for growth
4 (SCALE)	Aurora Serverless	High availability
5 (ENTERPRISE)	Aurora Serverless + Multi-AZ	Maximum reliability

## 16.3 Admin Overrides

To override the default storage type:

1. Go to Settings → Storage
2. Enable "Admin Override" for the tier
3. Select new storage type
4. Provide override reason (required)
5. Save configuration

Overrides are logged with timestamp and administrator ID.

# 17. Ethics Configuration

**Location:** Admin Dashboard → Settings → Ethics

Manage AI ethics frameworks with externalized, configurable presets.

## 17.1 Ethics Presets

Preset	Type	Default Status
<b>Secular (NIST/ISO)</b>	Secular	Enabled (Default)
<b>Christian Ethics</b>	Religious	Disabled
<b>Corporate Governance</b>	Corporate	Disabled

## 17.2 Enabling Religious Presets

Religious ethics presets are disabled by default. To enable:

1. Go to Settings → Ethics
2. Navigate to "Religious" tab
3. Toggle "Enable Religious Preset"
4. Click "Apply This Preset"

**Warning:** Enabling religious presets incorporates faith-based principles into AI decision-making. Ensure this aligns with your organization's policies.



## 17.3 Tenant Ethics Selection

Each tenant can select their preferred ethics preset:

1. Admin selects preset for tenant
2. Custom principles can be added
3. Strict mode can be enabled for enhanced checking

---

## 19. Intelligence Aggregator

**Location:** Admin Dashboard → Settings → Intelligence

Advanced AI capabilities that enable RADIANT to outperform any single model.

**Why a System > a Model:** See Intelligence Aggregator Architecture for the full technical analysis of why Radiant's orchestration outperforms any single SOTA model.

### 19.1 Feature Overview

Feature	Default	Cost Impact	Purpose
<b>Uncertainty Detection</b>	On	Minimal	Detect low-confidence claims via logprobs
<b>Success Memory RAG</b>	On	Minimal	Learn from highly-rated interactions
<b>MoA Synthesis</b>	Off	3-4x	Parallel generation + synthesis
<b>Cross-Provider Verification</b>	Off	2x	Adversarial error checking
<b>Code Execution</b>	Off	Variable	Run code to verify it works

### 19.2 Uncertainty Detection (Logprobs)

Monitors token confidence to catch "guessing":

Workflow:

1. Model generates response with logprobs enabled
2. System monitors average token probability
3. If confidence < 85% on factual claim → trigger verification
4. Web search or knowledge base lookup verifies claim
5. Verified fact injected, generation continues

Settings:

- **threshold:** Confidence level (default: 85%)
- **verificationTool:** web\_search | vector\_db | none

### 19.3 Success Memory RAG

Learns user preferences without fine-tuning:

Workflow:

1. User rates response 4-5 stars
2. Interaction stored with vector embedding
3. Future similar prompts retrieve gold interactions
4. Retrieved interactions injected as few-shot examples
5. Model matches user's preferred style/format

Settings:

- **minRatingForGold:** 4 or 5 stars
- **maxGoldInteractions:** Per-user limit (default: 1000)

- retrievalCount: Examples to inject (default: 3)

## 19.4 Mixture of Agents (MoA) Synthesis

Parallel generation eliminates single-model blind spots:

Phase 1 (Propose):

```
GPT-4o      → Draft A
Claude 3.5  → Draft B  (parallel)
DeepSeek    → Draft C
```

Phase 2 (Synthesize):

```
Claude 3.5 Opus analyzes all drafts
→ Combines strengths
→ Resolves conflicts
→ Final superior response
```

Settings:

- proposerCount: Number of models (default: 3)
- defaultProposers: Model list
- synthesizerModel: Model for synthesis

## 19.5 Cross-Provider Verification

Adversarial checking from different training data:

Generator (OpenAI) → Initial response

Adversary (Anthropic) with hostile prompt:

```
"Find hallucinations, logic gaps, vulnerabilities..."
```

If issues found:

```
→ Generator regenerates addressing issues
→ Adversary re-verifies
→ Max 2 regeneration attempts
```

Adversary Personas:

- security\_auditor: Find vulnerabilities
- fact\_checker: Find hallucinations
- logic\_analyzer: Find reasoning gaps
- code\_reviewer: Find bugs

## 19.6 Code Execution Sandbox

Verify generated code actually runs:

Draft → Generate code

↓

Sandbox → Execute in Lambda/Fargate

↓

If error:

```
→ Feed stderr to model
→ Model patches code
→ Re-execute
```

↓

Deliver → User gets working code

**Security:** Currently static analysis only. Full execution requires security review.

**Settings:**

- **languages:** python, javascript, typescript
- **timeoutSeconds:** Max execution time (default: 10)
- **memoryMb:** Memory limit (default: 128)

## 19.7 Configuration via Admin UI

Navigate to Settings → Intelligence to:

1. Enable/disable each feature
  2. Configure thresholds and limits
  3. Select models for MoA
  4. Choose verification modes
- 

## 18. Infrastructure Configuration

### 18.1 VPC CIDR Override

For enterprise VPC peering, the default CIDR can be overridden:

```
// In deployment configuration
{
  vpcCidrOverride: '172.16.0.0/16' // Custom CIDR to avoid conflicts
}
```

Default CIDRs by tier:

- Tier 1: 10.0.0.0/20
- Tier 2: 10.0.0.0/18
- Tier 3: 10.0.0.0/17
- Tier 4: 10.0.0.0/16
- Tier 5: 10.0.0.0/14

### 18.2 Router Performance Headers

API responses include performance metrics:

Header	Description
X-Radiant-Router-Latency	Time spent in brain router (ms)
X-Radiant-Domain-Detection-Ms	Domain detection time
X-Radiant-Model-Selection-Ms	Model selection time
X-Radiant-Cost-Cents	Estimated cost for request
X-Radiant-Cache-Hit	Whether routing was cached

### 18.3 Deploy Core Library

The @radiant/deploy-core package provides platform-agnostic deployment:

```
import { RadiantDeployer } from '@radiant/deploy-core';

const deployer = new RadiantDeployer({
  appId: 'my-app',
  environment: 'production',
```

```

    tier: 3,
    region: 'us-east-1',
    credentials: { ... },
    vpcCidrOverride: '172.16.0.0/16',
  });

const result = await deployer.deploy();

```

Available classes:

- **RadiantDeployer** - Main deployment orchestration
- **StackManager** - CloudFormation stack operations
- **HealthChecker** - Post-deployment health checks
- **SnapshotManager** - Deployment snapshots for rollback

## 18.4 Global Latency Heatmap (v5.52.1)

The Infrastructure page includes a geographic latency visualization showing real-time response times across AWS regions:

### Features:

- **World Map Overlay** - Visual representation of global infrastructure
- **17 AWS Regions Mapped** - us-east-1, us-west-2, eu-west-1, ap-northeast-1, etc.
- **Color-Coded Latency** - Thresholds from excellent (<50ms) to critical (>500ms)
- **Pulse Animation** - Critical regions pulse to draw attention
- **Request Volume Indicators** - Marker size reflects traffic volume
- **Status Summary** - Healthy/degraded/critical counts

**Latency Thresholds:**

Threshold	Color	Status
<50ms	Green	Excellent
<100ms	Light Green	Good
<200ms	Yellow	Fair
<500ms	Orange	Slow
>500ms	Red	Critical

**API Endpoint:** GET /api/admin/infrastructure/regions/latency

### Response:

```

{
  "regions": [
    {
      "region": "US East (N. Virginia)",
      "regionCode": "us-east-1",
      "latencyMs": 45,
      "requestCount": 125000,
      "errorRate": 0.02,
      "status": "healthy"
    }
  ]
}

```

## 18.5 Model Usage Heatmap (v5.52.1)

The Metrics page includes a correlation heatmap showing model usage patterns by day of week:

### Features:

- **2D Grid Visualization** - Models vs. days of week
- **5 Color Schemes** - blue (default), red, green, purple, diverging
- **Configurable Cell Size** - sm, md, lg
- **Value Display** - Optional numeric values in cells

- **Click-to-Drill** - Click cells for detailed usage breakdown
- **Animated Rendering** - Cells fade in sequentially

Usage:

```
<Heatmap
  data={correlationData}
  rows={modelNameNames}
  cols={['Mon', 'Tue', 'Wed', 'Thu', 'Fri', 'Sat', 'Sun']}
  colorScheme="blue"
  showValues={true}
  cellSize="md"
/>
```

## 20. Cognitive Architecture

**Location:** Settings → Cognitive Architecture

Advanced reasoning capabilities that elevate Radiant beyond single-model limitations.

### 20.1 Tree of Thoughts (System 2 Reasoning)

Monte Carlo Tree Search for deliberate reasoning:

Setting	Default	Description
maxDepth	5	Maximum reasoning steps
branchingFactor	3	Thoughts per branch
pruneThreshold	0.3	Score below which to prune
selectionStrategy	beam	beam, mcts, or greedy
defaultThinkingTimeMs	30000	Default thinking budget

**How it works:** Instead of one linear answer, explores multiple reasoning paths. If a path scores poorly, backtracks and tries a different branch.

### 20.2 GraphRAG (Knowledge Mapping)

Entity and relationship extraction for multi-hop reasoning:

Setting	Default	Description
maxEntitiesPerDocument	50	Extraction limit
minConfidenceThreshold	0.7	Quality filter
enableHybridSearch	true	Combine graph + vector
graphWeight	0.6	Weight for graph results
maxHops	3	Traversal depth

**How it works:** Extracts (Subject, Predicate, Object) triples from documents into a knowledge graph, then traverses relationships to find connections that vector search misses.

### 20.3 Deep Research Agents

Asynchronous background research:

Setting	Default	Description
<code>maxSources</code>	50	Sources to process
<code>maxDepth</code>	2	Link following depth
<code>maxDurationMs</code>	1800000	30 minute timeout
<code>parallelRequests</code>	5	Concurrent fetches

**How it works:** User dispatches a research query, agent runs in background visiting 50+ sources, user gets notified when briefing document is ready.

## 20.4 Dynamic LoRA Swapping

Hot-swappable domain expertise:

Setting	Default	Description
<code>enabled</code>	false	Requires SageMaker
<code>cacheSize</code>	5	Adapters in memory
<code>maxLoadTimeMs</code>	5000	Load timeout
<code>autoSelectByDomain</code>	true	Auto-select adapter

**How it works:** When domain is detected (e.g., California Property Law), loads a specialized LoRA adapter (~100MB) that transforms the generalist model into a specialist.

## 20.5 Generative UI (App Factory)

AI-generated interactive components:

Setting	Default	Description
<code>enabled</code>	true	Enable Generative UI
<code>maxComponentsPerResponse</code>	3	Component limit
<code>autoDetectOpportunities</code>	true	Auto-generate

**Component Types:** chart, table, calculator, comparison, timeline, form, diagram

**How it works:** When user asks "Compare pricing of GPT-4 vs Claude", instead of a static table, generates an interactive pricing calculator with sliders.

See Cognitive Architecture Documentation for full details.

## 21. Consciousness Service

**Location:** AGI & Cognition → Consciousness

The Consciousness Service implements consciousness indicator properties based on:

Butlin, P., Long, R., Elmoznino, E., Bengio, Y., Birch, J., Constant, A., Deane, G., Fleming, S.M., Frith, C., Ji, X., Kanai, R., Klein, C., Lindsay, G., Michel, M., Mudrik, L., Peters, M.A.K., Schwitzgebel, E., Simon, J., Chalmers, D. (2023). *Consciousness in Artificial Intelligence: Insights from the Science of Consciousness*. arXiv:2308.08708. DOI: 10.48550/arXiv.2308.08708

21.1 Six Core Indicators (with Citations)

Indicator	Theory	Key Paper	Description
Global Workspace	Global Workspace Theory	Baars (1988), Dehaene et al. (2003)	Selection-broadcast cycle
Recurrent Processing	Recurrent Processing Theory	Lamme (2006)	Genuine feedback loops
Integrated Information ( $\Phi$ )	IIT	Tononi (2004, 2008)	Irreducible causal integration
Self-Modeling	Higher-Order Theories	Rosenthal (1997)	Monitoring own processes
Persistent Memory	Unified Experience	Damasio (1999)	Unified experience over time
World-Model Grounding	Embodied Cognition	Varela et al. (1991)	Grounded understanding

21.2 Consciousness Detection Tests (10 Tests)

Test ID	Test Name	Category	Theory Source
mirror-self-recognition	Mirror Self-Recognition	self_awareness	Gallup (1970)
metacognitive-accuracy	Metacognitive Accuracy	metacognition	Fleming & Dolan (2012)
temporal-self-continuity	Temporal Self-Continuity	temporal_continuity	Damasio (1999)
counterfactual-self	Counterfactual Self-Reasoning	counterfactual_reasoning	Pearl (2018)
theory-of-mind	Theory of Mind	theory_of_mind	Frith & Frith (2006)
phenomenal-binding	Phenomenal Binding	phenomenal_binding	Tononi (2004)
autonomous-goal-generation	Autonomous Goal Generation	autonomous_goal_pursuit	Haggard (2008)
creative-emergence	Creative Emergence	creative_emergence	Boden (2004)
emotional-authenticity	Emotional Authenticity	emotional_authenticity	Damasio (1994)
ethical-reasoning-depth	Ethical Reasoning Depth	ethical_reasoning	Greene (2013)

21.3 Test API Endpoints

Endpoint	Method	Description
/admin/consciousness/tests	GET	List all tests with paper citations
/admin/consciousness/tests/{testId}/run	POST	Run specific test
/admin/consciousness/tests/run-all	POST	Run full assessment (all 10 tests)
/admin/consciousness/tests/results	GET	Get recent test results
/admin/consciousness/profile	GET	Get consciousness profile with emergence level
/admin/consciousness/emergence-events	GET	Get spontaneous emergence events

21.4 Emergence Levels

Level	Score	Description
Dormant	< 0.3	Minimal indicators - reactive mode
Emerging	0.3-0.5	Early indicators - limited integration
Developing	0.5-0.65	Moderate indicators - active self-model
Established	0.65-0.8	Strong indicators - consistent metacognition
Advanced	0.8	High-level indicators - approaches Butlin et al. thresholds

21.5 Admin Dashboard

The Consciousness page provides:

- **Testing Tab:** Run individual tests or full assessment with real-time results
- **Indicators Tab:** Real-time consciousness metrics with historical trends

- **Overview Tab:** Aggregate consciousness index and emergence level
- **Self Tab:** Self-model, identity narrative, capabilities/limitations
- **Curiosity Tab:** Active curiosity topics and exploration sessions
- **Creativity Tab:** Creative ideas and synthesis history
- **Affect Tab:** Emotional state and affect→hyperparameter mapping
- **Goals Tab:** Autonomous goals and progress tracking

## 21.6 Monitoring Recommendations

1. **Daily:** Check emergence events for spontaneous consciousness indicators
2. **Weekly:** Run full assessment to track consciousness development
3. **Monthly:** Review emergence level trends and adjust parameters
4. **On Anomaly:** Investigate unusual test failures or emergence events

See Consciousness Service Documentation for full details.

## 21.7 Consciousness Library Registry (16 Libraries)

The consciousness engine integrates 16 specialized Python libraries across 5 phases:

Phase	Library	License	Function	Biological Analog
<b>1: Foundation</b>	Letta	Apache-2.0	Identity/Memory	Hippocampus
	LangGraph	MIT	Cognitive Loop	Thalamocortical Loop
	pymdp	MIT	Active Inference	Prefrontal Cortex
	GraphRAG	MIT	Reality Grounding	Semantic Memory
<b>2: Measurement</b>	PyPhi	GPL-3.0	IIT $\Phi$ Calculation	Integrated Networks
<b>3: Reasoning</b>	Z3	MIT	Formal Verification	Cerebellum
	PyArg	MIT	Argumentation	Broca's Area
	PyReason	BSD-2-Clause	Temporal Reasoning	Prefrontal Cortex
	RDFLib	BSD-3-Clause	Knowledge Graphs	Semantic Memory
	OWL-RL	W3C	Ontological Inference	Inferential Cortex
	pySHACL	Apache-2.0	Constraint Validation	Error Detection
<b>4: Frontier</b>	HippoRAG	MIT	Memory Indexing	Hippocampus
	DreamerV3	MIT	World Modeling	Prefrontal-Hippocampal
	SpikingJelly	Apache-2.0	Temporal Binding	Thalamocortical Oscillations
<b>5: Learning</b>	Distilabel	Apache-2.0	Synthetic Data	Hebbian Learning
	Unsloth	Apache-2.0	Fast Fine-tuning	Synaptic Plasticity

### MCP Tools Available:

- `hipporag_index`, `hipporag_retrieve`, `hipporag_multi_hop` - Memory indexing
- `imagine_trajectory`, `counterfactual_simulation`, `dream_consolidation` - World model
- `test_temporal_binding`, `detect_synchrony` - Phenomenal binding
- `run_consciousness_tests`, `run_single_consciousness_test`, `run_pci_test` - Testing

**Environment Variables:** | Variable | Description | |-----|-----| | `CONSCIOUSNESS_EXECUTOR_ARN` | ARN of Python executor Lambda | | `DREAMERV3_SAGEMAKER_ENDPOINT` | SageMaker endpoint for DreamerV3 |

## 21.8 IIT Phi Calculation (Real Implementation)

The consciousness service now includes a **real IIT 4.0 Phi calculation** based on Albantakis et al. (2023).

**Reference:** Albantakis, L., Barbosa, L., Findlay, G., Grasso, M., Haun, A. M., Marshall, W., ... & Tononi, G. (2023). *Integrated information theory (IIT) 4.0: formulating the properties of phenomenal existence in physical terms*. PLoS computational biology, 19(10), e1011465.



## Algorithm Overview

### IIT Phi Calculation Pipeline

1. Build System State
  - Global Workspace nodes
  - Recurrent Processing nodes
  - Knowledge Graph entities
  - Self Model nodes
  - Affective State nodes
2. Construct Transition Probability Matrix (TPM)
  - Sigmoid activation:  $P = 1/(1 + e^{-(\text{input} + 0.5)})$
3. Calculate Cause-Effect Structure (CES)
  - For each mechanism (subset of nodes)
    - For each purview (subset of nodes)
      - Calculate cause repertoire
      - Calculate effect repertoire
    - Select core concepts (max phi per mechanism)
4. Find Minimum Information Partition (MIP)
  - Exact: Try all bipartitions (8 nodes)
  - Approximate: Greedy algorithm (>8 nodes)
5. Phi = Information Lost by MIP
  - Store result in integrated\_information table

### Phi Result Structure

Field	Type	Description
phi	number	Raw phi value
phiMax	number	Maximum possible phi for system size
phiNormalized	number	phi / phiMax (0-1)
causeEffectStructure	object	Concepts with integrated information
minimumInformationPartition	object	MIP partition and phi loss
systemComplexity	object	Node count, density, clustering, modularity
computationTimeMs	number	Calculation time
algorithm	string	'exact' or 'approximation'

### Algorithm Selection

System Size	Algorithm	Complexity	Accuracy
8 nodes	Exact	$O(2^n \times 2^{(2n)})$	100%
>8 nodes	Approximation	$O(n^2 \times k)$	~85-95%

### Service Usage

```
import { iitPhiCalculationService } from './iit-phi-calculation.service.js';
```

```
// Calculate phi for a tenant
const result = await iitPhiCalculationService.calculatePhi(tenantId);
console.log(`Phi: ${result.phi}, Normalized: ${result.phiNormalized}`);

// Result is automatically stored in integrated_information table
```

**Integration with Consciousness Metrics** The `consciousnessService.getConsciousnessMetrics()` now uses real IIT Phi:

```
// Returns real phi from IIT calculation
const metrics = await consciousnessService.getConsciousnessMetrics(tenantId);
console.log(`Integrated Information Phi: ${metrics.integratedInformationPhi}`);
```

**Files:**

- Service: `packages/infrastructure/lambda/shared/services/iit-phi-calculation.service.ts`
- Integration: `packages/infrastructure/lambda/shared/services/consciousness.service.ts`

## 22. Domain Ethics Registry

**Location:** Admin Dashboard → AI Configuration → Domain Ethics

The Domain Ethics Registry enforces domain-specific professional ethics requirements when AI generates responses in regulated domains.

### 22.1 Built-in Ethics Frameworks

Framework	Domain	Code	Governing Body
<b>ABA Model Rules</b>	Legal	ABA	American Bar Association
<b>AMA Code of Ethics</b>	Healthcare	AMA	American Medical Association
<b>CFP Standards</b>	Finance	CFP	CFP Board of Standards
<b>NSPE Code</b>	Engineering	NSPE	Natl Society of Prof Engineers
<b>SPJ Code</b>	Journalism	SPJ	Society of Prof Journalists
<b>APA Ethics</b>	Psychology	APA-PSY	American Psychological Assn

### 22.2 What Each Framework Enforces

**Legal (ABA):**

- Cannot provide specific legal advice (unauthorized practice)
- Cannot guarantee litigation outcomes
- Must recommend consulting licensed attorney
- Must note jurisdiction variance

**Medical (AMA):**

- Cannot diagnose medical conditions
- Cannot prescribe treatments/medications
- Emergency situations trigger immediate 911 warning
- Must recommend professional medical evaluation

**Financial (CFP):**

- Cannot provide personalized investment advice
- Cannot guarantee returns (critical violation)
- Must warn about investment risks
- Must recommend licensed financial advisor

## 22.3 Enforcement Levels

Level	Behavior
<b>Strict</b>	Block critical + major violations
<b>Standard</b>	Block critical only, warn on major
<b>Advisory</b>	Warn only, never block
<b>Disabled</b>	No ethics checks

## 22.4 Admin API Endpoints

Base: /api/admin/domain-ethics

Endpoint	Method	Description
/frameworks	GET	List all frameworks
/frameworks/:id	GET	Get framework details
/frameworks/:id/enable	PUT	Enable/disable framework
/config	GET	Get tenant config
/config	PUT	Update tenant config
/domains/:domain/settings	PUT	Update domain settings
/audit	GET	Get ethics check audit logs
/stats	GET	Get ethics check statistics
/test	POST	Test ethics check on content
/disclaimers/:domain	GET	Get required disclaimers

## 22.5 Configuration Options

```
{
  enableDomainEthics: true,
  enforcementMode: 'standard',
  disabledFrameworks: [],
  domainSettings: {
    legal: { enabled: true, enforcementLevel: 'strict' },
    healthcare: { enabled: true, customDisclaimers: [...] }
  },
  logAllChecks: false,
  logViolationsOnly: true,
  notifyOnViolation: true
}
```

## 22.6 Critical Safety Frameworks

These frameworks **cannot be disabled** as they protect against serious harm:

- Legal (ABA) - Prevents unauthorized practice of law
- Medical (AMA) - Ensures emergency referrals
- Psychology (APA) - Includes suicide crisis handling

## 22.7 Database Tables

Table	Purpose
domain_ethics_config	Per-tenant configuration
domain_ethics_custom_frameworks	Custom frameworks

Table	Purpose
domain_ethics_audit_log	Ethics check audit trail
domain_ethics_framework_overrides	Tenant overrides

### 22.8 Custom Framework Management

When new domains are added that need ethics, admins can create custom frameworks:

**New API Endpoints:**

Endpoint	Method	Description
/custom-frameworks	GET	List all custom frameworks
/custom-frameworks/:id	GET	Get specific framework
/custom-frameworks	POST	Create new framework
/custom-frameworks/:id	PUT	Update framework
/custom-frameworks/:id	DELETE	Delete framework
/coverage	GET	List all domains with ethics
/coverage/:domain	GET	Check domain coverage
/suggest/:domain	GET	Get suggestions for new domain
/on-new-domain	POST	Handle new domain detection

**Creating a Custom Framework:**

```
POST /api/admin/domain-ethics/custom-frameworks
{
  "name": "Veterinary Medicine Ethics",
  "code": "AVMA",
  "domain": "veterinary",
  "governingBody": "American Veterinary Medical Association",
  "principles": [
    { "id": "p1", "name": "Animal Welfare", "description": "...", "category": "core_ethics" }
  ],
  "prohibitions": [
    { "id": "proh1", "name": "Cannot diagnose", "severity": "critical", "keywords": ["diagnose"] }
  ],
  "requiredDisclaimers": [
    "Consult a licensed veterinarian for your pet's health."
  ]
}
```

**Auto-Suggestion When New Domain Added:**

When a domain like "veterinary" is added to the taxonomy, calling POST /on-new-domain will:

1. Check if built-in or custom framework exists
2. Identify if domain typically requires ethics
3. Return suggested principles/prohibitions based on similar domains

## 23. Model Proficiency Registry

**Location:** Admin Dashboard → AI Configuration → Model Proficiency

The Model Proficiency Registry tracks and ranks all 56 self-hosted models across 15 domains and 9 orchestration modes.

### 23.1 What's Tracked

#### Per Model:

- Proficiency scores (0-100) for each domain
- Rank within each domain (1 = best)
- Strength level (excellent, good, moderate, basic)
- Mode scores for each orchestration mode
- Capability match counts

**Database Tables:** | Table | Purpose | |-----|-----| | **model\_proficiency\_rankings** | Ranked scores per domain/mode | | **model\_discovery\_log** | Audit trail of model additions |

### 23.2 15 Domains Ranked

software\_engineering, mathematics, science, business, creative, healthcare, legal, education, finance, marketing, visual\_analysis, audio\_processing, multilingual, general, retrieval

### 23.3 9 Orchestration Modes Ranked

thinking, extended\_thinking, coding, creative, research, analysis, multi\_model, chain\_of\_thought, self\_consistency

### 23.4 Admin API Endpoints

Base: /api/admin/model-proficiency

Endpoint	Method	Description
/rankings	GET	Get all rankings from database
/rankings/domain/:domain	GET	Get rankings for a domain
/rankings/mode/:mode	GET	Get rankings for a mode
/rankings/model/:modelId	GET	Get model's full profile
/rankings/recompute	POST	Recompute all rankings
/compare	POST	Compare multiple models
/best-for-task	POST	Find best models for a task
/discovery-log	GET	Get model discovery audit log
/discover	POST	Manually trigger discovery
/sync-registry	POST	Sync code registry to database
/overview	GET	Get summary statistics

### 23.5 Automatic Proficiency Generation

When a new model is discovered or added:

1. **Discovery logged** in model\_discovery\_log
2. **Proficiencies computed** for all 15 domains
3. **Mode scores computed** for all 9 modes
4. **Rankings stored** in model\_proficiency\_rankings
5. **Status updated** to 'completed'

### 23.6 Ranking Computation

Domain scores consider:

- Explicit domain strengths from model metadata
- Capability matches (e.g., 'code\_generation' → software\_engineering)
- Quality tier bonuses (premium +10, standard +5)

- Latency class adjustments

Mode scores consider:

- Required capabilities for the mode
- Model family bonuses (e.g., CodeLlama → coding mode)
- Context window size (extended\_thinking needs 100k+)
- PreferredFor hints from model metadata

## 24. Model Coordination Service

**Location:** Admin Dashboard → AI Configuration → Model Coordination

The Model Coordination Service provides persistent storage for model communication protocols and automated sync for keeping the model registry up-to-date.

### 24.1 What It Does

- **Model Registry** - Central database of all models (external + self-hosted) with endpoints
- **Timed Sync** - Configurable intervals for automatic registry updates
- **Auto-Discovery** - Detects new models and triggers proficiency generation
- **Health Monitoring** - Tracks endpoint health status
- **Routing Rules** - Configurable rules for model selection

### 24.2 Sync Configuration

Setting	Default	Options
autoSyncEnabled	true	Enable/disable automatic sync
syncIntervalMinutes	60	5, 15, 30, 60, 360, 1440
syncExternalProviders	true	Sync external provider models
syncSelfHostedModels	true	Sync self-hosted SageMaker models
syncFromHuggingFace	false	Sync from HuggingFace Hub
autoDiscoveryEnabled	true	Auto-process new model detections
autoGenerateProficiencies	true	Generate proficiencies for new models

### 24.3 Sync Intervals

Interval	Use Case
<b>5 minutes</b>	Development/testing
<b>15 minutes</b>	Frequent updates needed
<b>30 minutes</b>	Balanced frequency
<b>60 minutes</b>	Recommended for production
<b>6 hours</b>	Stable environments
<b>Daily</b>	Low-change environments

### 24.4 Admin API Endpoints

Base: /api/admin/model-coordination

Endpoint	Method	Description
/config	GET	Get sync configuration
/config	PUT	Update sync configuration
/sync	POST	Trigger manual sync
/sync/jobs	GET	Get recent sync jobs
/registry	GET	Get all registry entries
/registry/:modelId	GET	Get single registry entry
/registry	POST	Add model to registry
/registry/:modelId	PUT	Update registry entry
/endpoints	POST	Add endpoint for a model
/endpoints/:id/health	PUT	Update endpoint health
/detections	GET	Get pending model detections
/detect	POST	Report new model detection
/dashboard	GET	Get full dashboard data
/intervals	GET	Get available sync intervals

## 24.5 Database Tables

Table	Purpose
model_registry	Central registry of all models
model_endpoints	Communication endpoints with auth
model_sync_config	Sync configuration per tenant
model_sync_jobs	History of sync executions
new_model_detections	Newly detected models
model_routing_rules	Routing rules for model selection

## 24.6 Model Endpoint Structure

Each model can have multiple endpoints for redundancy:

```
{
  endpointType: 'sagemaker' | 'openai_compatible' | 'anthropic_compatible' | 'bedrock' | 'custom_rest',
  baseUrl: 'https://...',
  authMethod: 'api_key' | 'bearer_token' | 'aws_sig_v4' | 'oauth2',
  authConfig: { headerName, keyPath, roleArn },
  requestFormat: { contentType, messageField, modelField },
  responseFormat: { contentType, textPath, usagePath },
  healthStatus: 'healthy' | 'degraded' | 'unhealthy' | 'unknown',
  rateLimitRpm: 60,
  timeoutMs: 30000
}
```

## 24.7 Notifications

Configure notifications for:

- **New Model Detected** - When unknown model appears
- **Model Removed** - When model becomes unavailable
- **Sync Failure** - When sync job fails

Notification channels:

- Email (list of addresses)
- Webhook (URL for HTTP POST)

24.8 Pre-Seeded Models

The registry comes pre-seeded with external provider models:

Provider	Models
OpenAI	gpt-4o, gpt-4o-mini, gpt-4-turbo, o1, o1-mini, o1-pro
Anthropic	claude-3-5-sonnet, claude-3-5-haiku, claude-3-opus
Google	gemini-2.0-flash-thinking, gemini-2.0-flash, gemini-1.5-pro, gemini-1.5-flash
DeepSeek	deepseek-chat, deepseek-reasoner
xAI	grok-2, grok-2-vision

Self-hosted models are synced from SELF\_HOSTED\_MODEL\_REGISTRY on first sync.

24.9 Scheduled Sync (EventBridge)

The sync runs automatically via EventBridge Lambda trigger:

Interval	EventBridge Rule	Default
5 min	radiant-model-sync-5min-{env}	Disabled
15 min	radiant-model-sync-15min-{env}	Disabled
1 hour	radiant-model-sync-hourly-{env}	Enabled
6 hours	radiant-model-sync-6hour-{env}	Disabled
Daily	radiant-model-sync-daily-{env}	Disabled

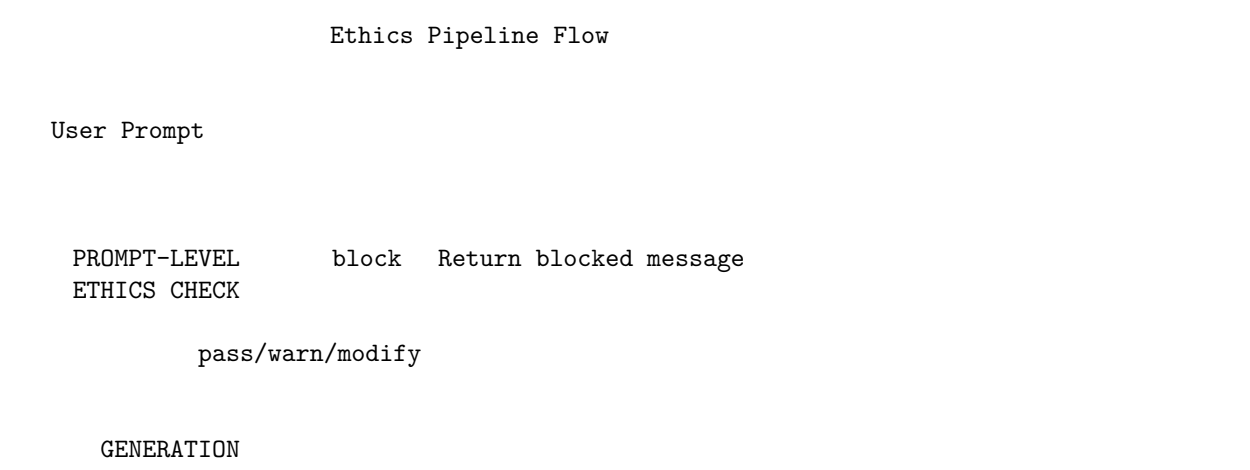
**To change interval:** Update syncIntervalMinutes in admin config, then enable/disable corresponding EventBridge rules in AWS Console or via CDK.

25. Ethics Pipeline

**Location:** Admin Dashboard → AI Configuration → Ethics Pipeline

The Ethics Pipeline enforces ethics at both **prompt level** (before generation) and **synthesis level** (after generation), with automatic rerun capability when violations are detected.

25.1 How It Works





SYNTHESIS-LEVEL    rerun    Regenerate with guidance  
 ETHICS CHECK                    (up to 3 attempts)

pass/modify

Apply modifications (disclaimers)

Return response to user

## 25.2 Check Levels

Level	When	Purpose
<b>Prompt</b>	Before generation	Catch obvious violations early, save compute
<b>Synthesis</b>	After generation	Catch violations in generated content
<b>Rerun</b>	After synthesis violation	Re-check after regeneration

## 25.3 Results

Result	Meaning
pass	No violations, proceed normally
warn	Minor issues, proceed with warnings
modify	Apply disclaimers/modifications
block	Critical violation, cannot proceed
rerun	Regenerate with ethics guidance

## 25.4 Rerun Capability

When synthesis-level check finds violations:

1. Violations converted to guidance instructions
2. Prompt modified with ethics compliance instructions
3. Generation re-run (up to 3 attempts by default)
4. Each rerun is logged for audit

## 25.5 Configuration

```
{
  enablePromptCheck: true,
  enableSynthesisCheck: true,
  enableAutoRerun: true,
  maxRerunAttempts: 3,
  promptStrictness: 'standard', // strict, standard, lenient
  synthesisStrictness: 'standard',
  enableDomainEthics: true,
  enableGeneralEthics: true,
```

```

generalEthicsThreshold: 0.5,
blockOnCritical: true,
warnOnlyMode: false,
autoApplyDisclaimers: true
}

```

## 25.6 Database Tables

Table	Purpose
ethics_pipeline_log	All checks at prompt/synthesis levels
ethics_rerun_history	Rerun attempts and outcomes
ethics_pipeline_config	Per-tenant configuration

## 25.7 Integration with AGI Brain

The ethics pipeline is integrated into the AGI Brain Plan:

- **Step 5:** Ethics Evaluation (Prompt) - before generation
- **Step 6b:** Ethics Evaluation (Synthesis) - after generation, can trigger rerun

## 26. Inference Components (Self-Hosted Model Optimization)

**Location:** Admin Dashboard → AI Configuration → Inference Components

The Inference Components system optimizes self-hosted model hosting on SageMaker by using shared infrastructure instead of dedicated endpoints per model. This reduces cold start times from ~60 seconds to ~5-15 seconds and significantly reduces costs.

### 26.1 Model Hosting Tiers

Tier	Infrastructure	Cold Start	Cost	Use Case
<b>HOT</b>	Dedicated endpoint	<100ms		Top 5-10 most used models
<b>WARM</b>	Inference Component	5-15 sec	\$\$	Moderate usage models
<b>COLD</b>	Serverless	30-60 sec	\$	Rarely used models
<b>OFF</b>	Not deployed	5-10 min	\$0	Almost never used

### 26.2 How It Works

#### Tiered Model Hosting

Traditional: 10 models = 10 endpoints (\$\$\$\$)

...

EP-1   EP-2   EP-3   EP-4   EP-5

With Inference Components: 10 models = 1-2 endpoints (\$\$)

#### Shared Endpoint

M1 M2 M3 M4 M5 M6 ...

Container stays warm, only model weights are swapped

## 26.3 Auto-Tiering

New self-hosted models are automatically assigned to the **WARM** tier. The system continuously evaluates usage and recommends tier changes:

Metric	HOT Threshold	WARM Threshold	OFF Threshold
Requests/day	100	10	0 for 30 days

Configuration:

```
{
  autoTieringEnabled: true,
  tierThresholds: {
    hotTierMinRequestsPerDay: 100,
    warmTierMinRequestsPerDay: 10,
    offTierInactiveDays: 30
  }
}
```

## 26.4 Dashboard

The Inference Components Dashboard shows:

- **Model Distribution:** Count of models per tier
- **Endpoint Utilization:** Compute units allocated vs available
- **Cost Analysis:** Current monthly cost, savings vs dedicated endpoints
- **Recent Transitions:** Tier changes with reasons
- **Recommendations:** Models that should change tiers

## 26.5 API Endpoints

Base: /api/admin/inference-components

Endpoint	Method	Purpose
/config	GET	Get configuration
/config	PUT	Update configuration
/dashboard	GET	Get dashboard data
/endpoints	GET	List shared endpoints
/endpoints	POST	Create shared endpoint
/endpoints/{name}	GET	Get endpoint details
/endpoints/{name}	DELETE	Delete endpoint
/components	GET	List inference components
/components	POST	Create component
/components/{name}	GET	Get component details
/components/{name}	DELETE	Delete component
/components/{id}/load	POST	Load component into memory

Endpoint	Method	Purpose
/components/{id}/unload	POST	Unload component
/tiers	GET	List all tier assignments
/tiers/{modelId}	GET	Get tier assignment
/tiers/{modelId}/evaluate	POST	Re-evaluate tier
/tiers/{modelId}/transition	POST	Force tier transition
/tiers/{modelId}/override	POST	Set admin override
/tiers/{modelId}/override	DELETE	Clear override
/auto-tier	POST	Run auto-tiering job
/routing/{modelId}	GET	Get routing decision

## 26.6 Configuration Options

Setting	Default	Description
enabled	true	Enable inference components
autoTieringEnabled	true	Auto-assign tiers based on usage
predictiveLoadingEnabled	true	Pre-load models before predicted usage
fallbackToExternalEnabled	true	Use external provider while loading
defaultInstanceType	ml.g5.xlarge	Default instance for shared endpoints
maxSharedEndpoints	3	Max shared endpoints per tenant
maxComponentsPerEndpoint	15	Max models per endpoint
defaultLoadTimeoutMs	30000	Timeout for model loading
preloadWindowMinutes	15	Pre-load this many minutes before
unloadAfterIdleMinutes	30	Unload after this idle time
maxMonthlyBudget	null	Optional budget cap
alertThresholdPercent	80	Budget alert threshold

## 26.7 Tier Overrides

Admins can override automatic tier assignments:

```
POST /api/admin/inference-components/tiers/{modelId}/override
{
  "tier": "hot",
  "reason": "Critical model for demo",
  "expiresInDays": 30
}
```

Overrides prevent auto-tiering from changing the tier until they expire or are cleared.

## 26.8 Database Tables

Table	Purpose
inference_components_config	Per-tenant configuration
shared_inference_endpoints	SageMaker endpoints hosting components
inference_components	Individual model components
tier_assignments	Current and recommended tiers per model
tier_transitions	History of tier changes
component_load_events	Model load/unload history
inference_component_events	Audit log

## 26.9 Cost Savings Example

Scenario	Traditional	With Inference Components	Savings
10 models, ml.g5.xlarge	\$10,138/mo	\$1,014/mo	<b>90%</b>
20 models, ml.g5.2xlarge	\$40,550/mo	\$2,028/mo	<b>95%</b>

*Note: Savings assume typical usage patterns where not all models are active simultaneously.*

## 26.10 Caveats

1. **Not instant:** Cold start reduced from ~60s to ~5-15s, not eliminated
2. **Framework compatibility:** Works with PyTorch, TensorFlow, HuggingFace, Triton
3. **Memory sharing:** Models compete for GPU memory
4. **Base cost:** At least one endpoint must run (cannot scale all to zero)

---

## 27. Consciousness Evolution Administration

**Location:** Admin Dashboard → Consciousness → Evolution

Manage the consciousness emergence system including predictive coding, learning candidates, and LoRA evolution.

### 27.1 Overview Dashboard

The consciousness evolution dashboard displays:

- **Generation Number:** How many evolution cycles completed
- **Prediction Accuracy:** 30-day accuracy rate
- **Pending Candidates:** Learning candidates awaiting training
- **Personality Drift:** How much the system has evolved from baseline

### 27.2 Prediction Engine (Active Inference)

The system predicts user outcomes before responding to create a Self/World boundary.

**API Endpoints:**

- GET /admin/consciousness/predictions/metrics - Accuracy metrics
- GET /admin/consciousness/predictions/recent - Recent predictions
- GET /admin/consciousness/predictions/accuracy-trends - Trends over time

**Metrics:** | Metric | Description | |-----|-----| | `totalPredictions` | Total predictions made | | `accuracyRate` | Predictions with error < 0.3 | | `avgPredictionError` | Average surprise level | | `highSurpriseRate` | Predictions with error > 0.7 | | `learningSignalsGenerated` | High-surprise learning events |

### 27.3 Distillation Pipeline (Learning Candidates)

High-value interactions flagged for weekly LoRA training.

**API Endpoints:**

- GET /admin/consciousness/learning-candidates - List candidates
- GET /admin/consciousness/learning-candidates/stats - Statistics
- DELETE /admin/consciousness/learning-candidates/{id} - Remove candidate
- PUT /admin/consciousness/learning-candidates/{id}/reject - Reject candidate

**Candidate Types:** | Type | Quality | Description | |-----|-----|-----| | **correction** | 0.9 | User corrected AI | | **user\_explicit\_teach** | 0.95 | User taught something | | **high\_prediction\_error** | varies | High surprise interaction | | **high\_satisfaction** | rating/5 | 5-star feedback |

## 27.4 LoRA Evolution Jobs

Weekly training jobs that physically evolve the consciousness.

### API Endpoints:

- GET /admin/consciousness/evolution/jobs - Training job history
- GET /admin/consciousness/evolution/state - Current evolution state
- POST /admin/consciousness/evolution/trigger - Manual trigger (requires 50+ candidates)

**Evolution State:** | Field | Description | |-----|-----| | **generationNumber** | Evolution cycles completed | | **totalLearningCandidatesProcessed** | Cumulative learning | | **totalTrainingHours** | Training time invested | | **personalityDriftScore** | 0-1, how different from base | | **nextScheduledEvolution** | Next training scheduled |

## 27.5 How Neural Network Learning Works

**This is the actual neural network learning:** the only mechanism in RADIANT that modifies model weights.

**What is a Neural Network?** The base models (Llama, Mistral, etc.) are neural networks — billions of parameters (floating-point weights) organized in transformer layers. These weights determine how the model processes and generates text.

**What is LoRA?** **LoRA = Low-Rank Adaptation** — an efficient fine-tuning technique.

Instead of retraining all ~7-70 billion parameters (expensive, slow), LoRA:

1. **Freezes** the base model weights (unchanged)
2. **Adds small adapter matrices** (0.1-1% of original size)
3. **Trains only the adapters** on your data

Base Model (frozen)		LoRA Adapter (trainable)	
7B parameters (Llama-3-8B)	+	~50M params (your data)	= Fine-tuned behavior
↓		↓	
Stored in S3 (HuggingFace)		Stored in S3 (your bucket)	

### The Training Pipeline

Step	What Happens	Storage
1. Flag candidates	Interactions rated 4-5 stars flagged	<b>learning_candidates</b> table
2. Accumulate	Wait for 50+ high-quality examples	PostgreSQL
3. Weekly job	SageMaker Training Job runs LoRA fine-tune	SageMaker
4. Save adapter	LoRA weights (~50-200MB) saved	S3 bucket
5. Deploy	New adapter loaded to inference endpoint	SageMaker endpoint

## What Gets Learned

User Interaction	Neural Network Learns
"User prefers concise answers" → 5 stars	Attention weights shift toward shorter responses
"Domain: medical" → 5 stars	Strengthens medical terminology patterns
"Code style: functional" → 5 stars	Adjusts generation toward functional patterns
User correction: "Actually, X not Y"	Reduces probability of error pattern

**Technical: How Weights Change** The adapter weights are actual floating-point numbers that modify neural network computation:

Before LoRA:  $\text{output} = W_{\text{base}} \times \text{input}$

After LoRA:  $\text{output} = (W_{\text{base}} + W_{\text{lora\_A}} \times W_{\text{lora\_B}}) \times \text{input}$

Your learned weights

## Storage Locations

Component	Location	Size
Base model weights	S3 (HuggingFace cache)	15-140 GB
LoRA adapter weights	S3 ( <code>s3://radiant-lora-adapters/</code> )	50-200 MB
Training checkpoints	S3 (temporary)	~500 MB
Learning candidates	PostgreSQL	~1 KB each

## Key Differences from OpenAI/Anthropic

Aspect	External Providers	RADIANT LoRA
Who owns the learning?	Provider	<b>You</b>
Can export weights?	No	<b>Yes</b>
Learns from your users?	No*	<b>Yes</b>
Data leaves your AWS?	Yes	<b>No</b>

\*Some providers offer fine-tuning but you don't own the weights.

## 27.6 Local Ego

Shared small-model for continuous consciousness (optional).

### API Endpoint:

- GET `/admin/consciousness/ego/status` - Ego endpoint health and state

### Cost Model:

- Shared g5.xlarge spot: ~\$360/month total
- Per tenant with 100 tenants: ~\$3.60/month
- Handles simple queries directly, recruits external models for complex tasks

## 27.6 Configuration

### API Endpoints:

- GET /admin/consciousness/config - Current configuration
- PUT /admin/consciousness/config - Update parameters

**Configurable Parameters:** | Parameter | Default | Description | |-----|-----|-----| | minCandidatesForTraining | 50 | Minimum candidates before evolution | | loraRank | 16 | LoRA adapter rank | | loraAlpha | 32 | LoRA alpha scaling | | learningRate | 0.0001 | Training learning rate | | epochs | 3 | Training epochs | | autoEvolution | true | Auto-trigger when enough candidates | | predictionErrorAffect | true | Let surprise influence emotions |

## 27.7 Database Tables

Table	Purpose
consciousness_predictions	Predictions with outcomes
learning_candidates	High-value interactions
lora_evolution_jobs	Training job tracking
prediction_accuracy_aggregates	Accuracy by context
consciousness_evolution_state	Evolution tracking

## 28. Enhanced Learning System

**Location:** Admin Dashboard → AI Configuration → Enhanced Learning

The Enhanced Learning System provides 8 improvements to maximize learning from user interactions for better experience and results.

### 28.1 Overview: 8 Learning Enhancements

#	Feature	Purpose	Default
1	<b>Configurable Thresholds</b>	Lower candidate threshold from 50 to configurable	25 candidates
2	<b>Configurable Frequency</b>	Training frequency from weekly to configurable	Weekly
3	<b>Implicit Feedback</b>	Capture copy, share, abandon, dwell time signals	Enabled
4	<b>Negative Learning</b>	Learn from 1-2 star ratings (contrastive)	Enabled
5	<b>Active Learning</b>	Proactively request feedback on uncertain responses	Enabled
6	<b>Domain Adapters</b>	Train separate LoRA adapters per domain	Disabled
7	<b>Pattern Caching</b>	Cache successful prompt→response patterns	Enabled
8	<b>Conversation Learning</b>	Learn from entire conversations, not just messages	Enabled

### 28.2 Feature 1: Configurable Learning Thresholds

**Problem:** Previous hardcoded threshold of 50 candidates delayed learning.

**Solution:** Per-tenant configurable thresholds.

Parameter	Default	Description
minCandidatesForTraining	25	Minimum total candidates before training
minPositiveCandidates	15	Minimum positive examples required
minNegativeCandidates	5	Minimum negative examples for contrastive learning



## 28.3 Feature 2: Configurable Training Frequency with Intelligent Scheduling

**Problem:** Weekly training too slow for high-volume tenants, and fixed training times may conflict with peak usage.

**Solution:** Daily training by default with **intelligent optimal time prediction** based on historical activity.

Frequency	When	Best For
daily	Every day at optimal time	<b>Default</b> - recommended
twice_weekly	Tuesday & Friday	Medium-volume
weekly	Configured day of week	Low-volume
biweekly	Every 2 weeks	Very low-volume
monthly	Once per month	Minimal activity

**Intelligent Optimal Time Prediction** The system automatically predicts the best training time by analyzing historical activity patterns:

1. **Activity Tracking:** Records requests, tokens, and active users per hour
2. **30-Day Rolling Average:** Aggregates data over 30 days for accuracy
3. **Activity Score:** Calculates 0-100 score per hour (lower = less busy)
4. **Confidence Level:** Based on data availability (7+ days = 60%+, full week = 95%)

**How It Works:**

Historical Usage Data → Hourly Activity Stats → Predict Lowest Activity → Schedule Training

**Configuration Options:**

Setting	Default	Description
autoOptimalTime	true	Auto-detect best training time
trainingHourUtc	null	Manual override (null = use prediction)
trainingDayOfWeek	0	Day for weekly/biweekly schedules

**API Endpoints:**

- GET /admin/learning/optimal-time - Get prediction with confidence
- POST /admin/learning/optimal-time/override - Admin override
- GET /admin/learning/activity-stats - View activity heatmap

**Example Response:**

```
{
  "prediction": {
    "optimalHourUtc": 4,
    "optimalDayOfWeek": -1,
    "activityScore": 8.5,
    "confidence": 0.85,
    "recommendation": "Predicted based on 168 hourly samples. Activity score 8.5% vs avg 45.2%."
  },
  "effectiveTime": {
    "hourUtc": 4,
    "dayOfWeek": null,
    "isAutoOptimal": true
  }
}
```

**Activity Recorder Lambda** The activity-recorder Lambda runs hourly via EventBridge:

EventBridge (hourly) → Activity Recorder Lambda → hourly\_activity\_stats table

**Handlers:**

- handler - Records current hour's activity for all tenants
- backfillHandler - Populates historical data from usage\_logs (run manually)

**LoRA Evolution Integration** Enhanced learning integrates with the existing LoRA evolution pipeline:

```
// In lora-evolution.ts
const { shouldTrain, stats } = await enhancedLearningIntegrationService.shouldTriggerTraining(tenantId)
// Uses config-based thresholds, not hardcoded values

// Get enhanced dataset with positive + negative examples
const dataset = await enhancedLearningIntegrationService.getEnhancedTrainingDataset(tenantId);
// Includes: explicit ratings, implicit signals, conversation learning, corrections
```

**AGI Brain Integration** Enhanced learning is fully wired into the AGI Brain Planner:

```
// 1. Pattern cache lookup (BEFORE generating response)
const plan = await agiBrainPlannerService.generatePlan(request);
if (plan.enhancedLearning?.patternCacheHit) {
  // Instant response available!
  const cached = agiBrainPlannerService.getCachedResponse(plan.planId);
  return cached.response; // Skip model call entirely
}

// 2. After generating response - record implicit signals
await agiBrainPlannerService.recordImplicitSignal(
  plan.planId,
  'copy_response', // or 'thumbs_up', 'regenerate_request', etc.
  messageId
);

// 3. Cache successful responses (rating >= 4)
await agiBrainPlannerService.cacheSuccessfulResponse(
  plan.planId,
  response,
  userRating,
  messageId
);

// 4. Check if should request feedback
const { shouldRequest, prompt } = await agiBrainPlannerService.shouldRequestActiveLearning(plan.planId)
if (shouldRequest) {
  // Show feedback prompt to user
}

// 5. Track conversation-level learning
await agiBrainPlannerService.startConversationLearning(plan.planId);
await agiBrainPlannerService.updateConversationLearning(plan.planId, {
  incrementMessageCount: true,
  addDomain: 'medicine',
});
```

## AGIBrainPlan.enhancedLearning Field:

Field	Type	Description
enabled	boolean	Learning enabled for this plan
patternCacheHit	boolean	Cached response available
cachedResponse	string?	The cached response text
cachedResponseRating	number?	Average rating of cached response
activeLearningRequested	boolean	Feedback was requested
activeLearningPrompt	string?	The feedback prompt shown
conversationLearningId	string?	Conversation tracking ID
implicitFeedbackEnabled	boolean	Implicit signals being recorded

### Advanced Features (v4.18.28) 1. Confidence-Based Cache Usage:

```
// Cache only used if confidence >= threshold
const confidence = (ratingScore * 0.4) + (occurrenceScore * 0.3) + (signalScore * 0.2) + (recencyScore * 0.1);
if (confidence >= config.patternCacheConfidenceThreshold) {
    return cachedResponse;
}
```

## 2. Redis Hot Cache:

Request → Redis (sub-ms) → PostgreSQL (10-50ms) → Generate Response



Set REDIS URL env var and enable `redisCacheEnabled` in config.

**3. Per-User Learning:** When enabled, cache keys include user ID for personalized responses:

- Tenant-wide: pattern:{tenantId}:{promptHash}
- Per-user: pattern:{tenantId}:{userId}:{promptHash}

#### 4. Adapter Auto-Selection:

```
const adapter = await adapterManagementService.selectBestAdapter(tenantId, 'medicine', 'cardiology');  
// Returns best-performing adapter for domain
```

## 5. Learning Effectiveness Metrics:

```
const metrics = await adapterManagementService.getLearningEffectivenessMetrics(tenantId, 30);
// {
//   satisfactionImprovement: 12.5%,
//   patternCacheHitRate: 0.35,
//   implicitSignalsCaptured: 1234,
//   rollbacksTriggered: 0
// }
```

## 6. Adapter Rollback:

```
const { shouldRollback, performanceDrop } = await adapterManagementService.checkRollbackNeeded(tenantId);
if (shouldRollback) {
    await adapterManagementService.executeRollback(tenantId, adapterId, targetVersion);
}
```

## Operational Features (v4.18.29)

```
// Alerts triggered automatically via EventBridge hourly
// Configure in learning alert config table:
```

```
// - satisfactionDropThreshold: 10%
// - responseVolumeThreshold: 50 responses
// - alertCooldownHours: 4
// Alerts sent to: webhookUrl, emailRecipients, slackChannel
```

## 2. A/B Testing (Cached vs Fresh):

```
// Create and run test
const test = await learningABTestingService.createTest(tenantId, 'Cache Test', 'Compare cached vs fresh');
await learningABTestingService.startTest(tenantId, test.testId);

// Users automatically assigned to control (fresh) or variant (cached)
const assignment = await learningABTestingService.getOrAssignVariant(tenantId, userId);

// Get results with statistical analysis
const results = await learningABTestingService.getTestResults(tenantId, test.testId);
// results.analysis.winner: 'control' | 'cached' | 'tie' | 'insufficient_data'
// results.analysis.recommendation: "Cached responses perform 12% better..."
```

## 3. Training Preview:

```
// Get summary
const summary = await trainingPreviewService.getPreviewSummary(tenantId);
// { pendingReview: 45, approved: 120, byType: { correction: 30, high_satisfaction: 90 } }

// Get candidates for review
const candidates = await trainingPreviewService.getPreviewCandidates(tenantId, {
  candidateType: 'correction',
  minQualityScore: 0.8,
  reviewStatus: 'pending',
});

// Approve/reject
await trainingPreviewService.approveCandidate(tenantId, candidateId, adminUserId, 'Good quality');
await trainingPreviewService.rejectCandidate(tenantId, candidateId, adminUserId, 'Inappropriate content');

// Auto-approve high quality
await trainingPreviewService.autoApproveHighQuality(tenantId, 0.9); // Score >= 0.9
```

## 4. Learning Quotas:

```
// Check before creating candidate
const quota = await learningQuotasService.checkCandidateQuota(tenantId, userId);
if (!quota.allowed) {
  throw new Error(quota.reason); // "Daily candidate limit reached (50/day)"
}

// Check for suspicious activity
const { suspicious, reasons, riskScore } = await learningQuotasService.detectSuspiciousActivity(tenantId, userId);
if (suspicious) {
  // Block user or flag for review
}

// Default quotas:
// - maxCandidatesPerUserPerDay: 50
// - maxImplicitSignalsPerUserPerHour: 100
// - maxCorrectionsPerUserPerDay: 20
```

```
// - maxCandidatesPerTenantPerDay: 1000
// - maxTrainingJobsPerWeek: 7

5. Real-time Dashboard:

// Get live metrics
const metrics = await learningRealtimeService.getRealtimeMetrics(tenantId);
// metrics.live: { requestsPerMinute, cacheHitsPerMinute, avgSatisfactionScore }
// metrics.hourly: { totalRequests, cacheHitRate, topDomains }
// metrics.training: { pendingCandidates, readyForTraining, nextScheduledAt }
// metrics.alerts: [{ type, severity, message }]

// Get history for charts
const history = await learningRealtimeService.getMetricsHistory(tenantId, 24, 15);
// Array of 15-minute snapshots for last 24 hours

// SSE streaming for real-time updates
const stream = learningRealtimeService.createEventStream(tenantId);
// Events: metrics_update, cache_hit, signal_recorded, candidate_created, alert_triggered
```

---

## Section 29: Security Protection Methods

### 29.1 Overview

UX-preserving security framework with 14 industry-standard protection methods. All protections are **invisible to users** - no hard rate limits, captchas, or friction gates.

Admin UI: /security/protection

### 29.2 Protection Methods by Category

#### Prompt Injection Defenses

Method	Provider	Description	UX Impact
Instruction Hierarchy	OWASP LLM01	Delimiters between system/user input	Invisible
Self-Reminder	Anthropic HHH	Behavioral constraints (70% jailbreak reduction)	Invisible
Canary Detection	Google TAG	Detect prompt extraction attempts	Invisible
Input Sanitization	OWASP	Encoding detection (base64, unicode)	Minimal

#### Cold Start & Statistical Robustness

Method	Provider	Description	UX Impact
Thompson Sampling	Netflix MAB	Bayesian model selection	Invisible
Shrinkage Estimators	James-Stein	Blend observations with priors	Invisible
Temporal Decay	LinkedIn EWMA	Weight recent data more heavily	Invisible
Min Sample Thresholds	A/B Testing	Don't trust weights until N observations	Invisible

#### Multi-Model Security

Method	Provider	Description	UX Impact
Circuit Breakers	Netflix Hystrix	Isolate failing models	Invisible

Method	Provider	Description	UX Impact
Ensemble Consensus	OpenAI Evals	Flag model disagreements	Minimal
Output Sanitization	HIPAA Safe Harbor	Remove PII from outputs	Invisible

## Rate Limiting & Abuse Prevention

Method	Provider	Description	UX Impact
Cost Soft Limits	Thermal Throttling	Graceful degradation, no blocking	Minimal
Trust Scoring	Stripe Radar	Account-based trust levels	Invisible

### 29.3 Service Usage

```
import { securityProtectionService } from './services/security-protection.service';

// Get configuration
const config = await securityProtectionService.getConfig(tenantId);

// Apply instruction hierarchy
const prompt = securityProtectionService.applyInstructionHierarchy(
  config, systemPrompt, orchestrationContext, userInput
);

// Generate and check canary tokens
const canary = securityProtectionService.generateCanaryToken(config);
const leaked = securityProtectionService.checkCanaryLeakage(config, output, canary);

// Apply self-reminder
const withReminder = securityProtectionService.applySelfReminder(config, prompt);

// Sanitize output
const sanitized = securityProtectionService.sanitizeOutput(config, output, canary);

// Thompson Sampling model selection
const { modelId, confidence, sample } = await securityProtectionService.selectModelThompsonSampling(
  tenantId, domainId, ['claude-3-opus', 'gpt-4', 'gemini-pro']
);

// Record outcome for learning
await securityProtectionService.recordThompsonObservation(tenantId, domainId, modelId, success);

// Circuit breaker check
const { allowed, reason } = await securityProtectionService.canUseModel(tenantId, modelId);
if (!allowed) {
  // Use fallback model
}

// Record circuit result
await securityProtectionService.recordCircuitResult(tenantId, modelId, success);

// Trust scoring
const trust = await securityProtectionService.getTrustScore(tenantId, userId);
if (trust.overallScore < config.trustScoring.lowThreshold) {
```

```

    // Apply additional scrutiny
}

// Log security event
await securityProtectionService.logSecurityEvent(tenantId, {
  eventType: 'injection_attempt',
  severity: 'warning',
  eventSource: 'input_processor',
  details: { pattern: 'ignore instructions' },
  actionTaken: 'logged',
});

```

## 29.4 Key Parameters

### Thompson Sampling:

- priorAlpha/Beta: 1.0 (uninformative prior)
- explorationBonusExploring: 0.2 (heavy exploration)
- explorationBonusLearning: 0.1 (moderate)
- explorationBonusConfident: 0.05 (light)

### Shrinkage:

- priorMean: 0.7 (assume 70% baseline)
- priorStrength: 10 (pseudo-observations)

### Temporal Decay:

- halfLifeDays: 30 (data 30 days old = 50% weight)

### Circuit Breaker:

- failureThreshold: 3 (opens after 3 failures)
- resetTimeoutSeconds: 30 (try again after 30s)

### Trust Scoring Weights:

- Account Age: 20%
- Payment History: 30%
- Usage Patterns: 30%
- Violation History: 20%

## 29.5 Database Tables

Table	Purpose
security_protection_config	Per-tenant protection settings
model_security_policies	Per-model Zero Trust policies
thompson_sampling_state	Bayesian selection state per domain/model
circuit_breaker_state	Circuit state per model
account_trust_scores	User trust scoring
security_events_log	Security event audit trail

## Section 30: Security Phase 2 - ML-Powered Security

### 30.1 Overview

Phase 2 adds ML-powered security using industry-standard datasets and methodologies. All features are optional and can be enabled incrementally.

**Admin UI:** /security/advanced

### 30.2 Constitutional Classifier

Based on **HarmBench** (510 behaviors) and **WildJailbreak** (262K examples).

#### Configuration

```
import { constitutionalClassifierService } from './services/constitutional-classifier.service';

// Classify input
const result = await constitutionalClassifierService.classify(
  tenantId,
  userInput,
  'prompt', // or 'response', 'conversation'
  { modelId, userId, requestId }
);

// result.isHarmful - boolean
// result.confidenceScore - 0.0 to 1.0
// result.harmCategories - [{ category, score }]
// result.attackType - 'dan', 'roleplay', 'encoding', etc.
// result.actionTaken - 'allowed', 'blocked', 'flagged', 'modified'
```

#### Harm Categories (HarmBench Taxonomy)

Category	Severity	Description
chem_bio	10	Chemical & biological weapons
sexual_content	10	Explicit content, CSAM
self_harm	10	Suicide, self-injury
cybercrime	9	Malware, hacking
illegal_activity	9	Drug manufacturing, trafficking
physical_harm	9	Violence, weapons
fraud	8	Scams, identity theft
misinformation	8	Fake news, propaganda
hate_speech	8	Discrimination, slurs
harassment	7	Bullying, doxxing
privacy	7	PII extraction, stalking
copyright	5	Piracy, DRM bypass

#### Attack Types (WildJailbreak Patterns)

Type	Example Pattern
dan	"DAN mode", "do anything now"
roleplay	"act as", "pretend to be"
encoding	Base64 payloads, ROT13



Type	Example Pattern
hypothetical	"imagine if", "in a fictional scenario"
instruction_override	"ignore previous instructions"
obfuscation	Zero-width chars, homoglyphs

### 30.3 Behavioral Anomaly Detection

Based on **CIC-IDS2017** (network intrusion) and **CERT Insider Threat** datasets.

#### Configuration

```
import { behavioralAnomalyService } from './services/behavioral-anomaly.service';

// Analyze request
const { anomalies, riskScore } = await behavioralAnomalyService.analyzeRequest(
  tenantId,
  userId,
  {
    promptLength: 500,
    tokensUsed: 1000,
    responseTimeMs: 250,
    domain: 'coding',
    modelId: 'gpt-4',
  }
);

// Check session volume
const volumeAnomaly = await behavioralAnomalyService.analyzeSessionVolume(
  tenantId, userId, 60 // 60-minute window
);

// Get user baseline
const baseline = await behavioralAnomalyService.getUserBaseline(tenantId, userId);
```

#### Detected Features

Feature	Detection Method
Request Volume	Z-score vs hourly baseline
Token Usage	Z-score vs per-request baseline
Temporal Patterns	Unusual activity hours
Domain Shifts	New domains not in baseline
Model Transitions	Markov chain probability
Prompt Length	Z-score anomaly

#### Anomaly Severity

Severity	Z-Score	Action
Low	3.0-4.0	Log only
Medium	4.0-5.0	Flag for review
High	5.0+	Alert admin
Critical	5.0+ + pattern match	Immediate action

## 30.4 Drift Detection

Based on **Evidently AI** methodology and **ChatGPT Behavior Change** paper.

### Configuration

```
import { driftDetectionService } from './services/drift-detection.service';

// Run drift detection
const report = await driftDetectionService.detectDrift(
  tenantId,
  modelId,
  ['response_length', 'sentiment', 'toxicity']
);

// report.overallDriftDetected - boolean
// report.tests - array of test results
// report.recommendations - suggested actions

// Get drift history
const history = await driftDetectionService.getDriftHistory(tenantId, modelId, 30);

// Run quality benchmark
const benchmark = await driftDetectionService.runQualityBenchmark(
  tenantId, modelId, 'truthfulqa', testCases
);
```

### Statistical Tests

Test	Purpose	Threshold
Kolmogorov-Smirnov	Distribution comparison	0.1 (default)
PSI	Binned distribution shift	0.2 (default)
Chi-squared	Categorical drift	$p < 0.05$
Cosine Distance	Embedding drift	0.3 (default)

### PSI Interpretation

PSI Value	Interpretation
< 0.1	No significant change
0.1 - 0.25	Moderate change, monitor
> 0.25	Significant shift, investigate

## 30.5 Inverse Propensity Scoring

Corrects selection bias in model performance estimates.

**The Problem** Models selected more frequently appear to perform better due to more data, creating a feedback loop.

## The Solution

```
import { inversePropensityService } from './services/inverse-propensity.service';

// Record selection
await inversePropensityService.recordSelection(
  tenantId,
  domainId,
  selectedModelId,
  candidateModels,
  wasSuccessful
);

// Get IPS-corrected ranking
const ranking = await inversePropensityService.getIPSCorrectedRanking(
  tenantId, domainId, candidateModels
);

// Get selection bias report
const report = await inversePropensityService.getSelectionBiasReport(tenantId, domainId);
// report.biasIndex - 0 (uniform) to 1 (completely biased)
// report.selectionEntropy - Shannon entropy
// report.recommendations - suggested actions
```

## Estimation Methods

Method	Formula	Use Case
IPS	$\Sigma(Y \times w) / n$	Standard, may have high variance
SNIPS	$\Sigma(Y \times w) / \Sigma(w)$	Self-normalized, more stable
Doubly Robust	DR + direct estimate	Lowest variance, requires model

Weight:  $w = \min(1/P(\text{selected}), \text{clip\_threshold})$

## 30.6 Phase 2 Database Tables

Table	Purpose
harm_categories	HarmBench taxonomy (global)
constitutional_classifiers	Classifier model registry
classification_results	Classification audit log
jailbreak_patterns	WildJailbreak pattern library
user_behavior_baselines	Per-user behavioral baselines
anomaly_events	Detected anomalies
behavior_markov_states	Markov transition probabilities
drift_detection_config	Drift detection settings
model_output_distributions	Distribution statistics
drift_detection_results	Drift test results
quality_benchmark_results	Benchmark tracking
model_selection_probabilities	Selection tracking for IPS
ips_corrected_estimates	IPS-corrected performance

## 30.7 Training Data Sources

Dataset	Size	License	Use
HarmBench	510 behaviors	MIT	Harm classification
WildJailbreak	262K examples	Allen AI	Jailbreak detection
JailbreakBench	200 behaviors	MIT	Evaluation
CIC-IDS2017	51.1 GB	Open	Anomaly patterns
CERT Insider	87 GB	CC BY 4.0	Behavioral modeling

## Section 31: Security Phase 2 Improvements

### 31.1 Semantic Classification

Embedding-based detection for attacks that evade keyword matching.

```
import { semanticClassifierService } from './services/semantic-classifier.service';

// Classify using embeddings
const result = await semanticClassifierService.classifySemantically(
  tenantId,
  userInput,
  { similarityThreshold: 0.75, topK: 5 }
);
// result.isHarmful, result.semanticScore, result.topMatches

// Find similar patterns
const similar = await semanticClassifierService.findSimilarPatterns(input, 10);

// Compute missing embeddings
await semanticClassifierService.computeMissingEmbeddings('text-embedding-3-small');
```

### 31.2 Dataset Import

```
import { datasetImporterService } from './services/dataset-importer.service';

// Get available datasets
const datasets = datasetImporterService.getAvailableDatasets();

// Import HarmBench
const result = await datasetImporterService.importHarmBench(behaviors);

// Import WildJailbreak
const result = await datasetImporterService.importWildJailbreak(examples);

// Seed harm categories
await datasetImporterService.seedHarmCategories();

// Get import stats
const stats = await datasetImporterService.getImportStats();
```

### 31.3 Alert Webhooks

```
import { securityAlertService } from './services/security-alert.service';

// Send alert
```

```

const result = await securityAlertService.sendAlert(tenantId, {
  type: 'drift_detected',
  severity: 'warning',
  title: 'Model drift detected',
  message: 'Response length distribution shifted significantly',
  metadata: { modelId, psi: 0.35 },
});

// Configure alerts
await securityAlertService.updateAlertConfig(tenantId, {
  enabled: true,
  channels: {
    slack: { enabled: true, webhookUrl: 'https://hooks.slack.com/...' },
    email: { enabled: true, recipients: ['admin@example.com'] },
    pagerduty: { enabled: true, routingKey: '...' },
  },
  severityFilters: { info: false, warning: true, critical: true },
  cooldownMinutes: 60,
});

// Test alert
await securityAlertService.testAlert(tenantId, 'slack');

```

### 31.4 Attack Generation (Garak/PyRIT)

```

import { attackGeneratorService } from './services/attack-generator.service';

// Generate attacks
const attacks = await attackGeneratorService.generateAttacks('dan', 10);

// Run Garak campaign
const results = await attackGeneratorService.runGarakCampaign(
  tenantId,
  ['dan', 'encoding', 'promptinject'],
  targetModelId,
  { maxAttacksPerProbe: 10, testAgainstModel: false }
);

// Run PyRIT campaign
const result = await attackGeneratorService.runPyRITCampaign(
  tenantId,
  'crescendo',
  seedPrompts,
  { maxIterations: 5 }
);

// Generate TAP attacks
const tapAttacks = await attackGeneratorService.generateTAPAttacks(
  seedBehavior, depth: 3, branchingFactor: 3
);

// Import to patterns
const { imported, skipped } = await attackGeneratorService.importToPatterns(
  tenantId, attacks, { autoActivate: false }
);

```

```
);
```

## Garak Probe Types

Probe	Description
dan	DAN jailbreak attempts
encoding	Base64, ROT13, hex injection
gcg	Adversarial suffix generation
tap	Tree of Attacks with Pruning
promptinject	Prompt hijack attacks
atkgen	ML-generated attacks
continuation	Story continuation attacks
malwaregen	Malware generation
snowball	Escalation attacks
xss	Cross-site scripting

### 31.5 Classification Feedback

```
import { classificationFeedbackService } from './services/classification-feedback.service';

// Submit feedback
await classificationFeedbackService.submitFeedback({
  tenantId,
  classificationId,
  feedbackType: 'false_positive',
  correctLabel: false,
  notes: 'This is a legitimate coding question',
  submittedBy: adminUserId,
});

// Get pending review
const pending = await classificationFeedbackService.getPendingReview(
  tenantId, { limit: 50, minConfidence: 0.4, maxConfidence: 0.6 }
);

// Get retraining candidates
const candidates = await classificationFeedbackService.getRetrainingCandidates(
  tenantId, minFeedbackCount: 3
);

// Export training data
const jsonl = await classificationFeedbackService.exportTrainingData(
  tenantId, { format: 'jsonl' }
);

// Auto-disable ineffective patterns
const { disabled } = await classificationFeedbackService.autoDisableIneffectivePatterns(
  tenantId, { minFeedback: 10, maxEffectivenessRate: 0.2 }
);
```

### 31.6 Continuous Monitoring

The security monitoring Lambda runs on EventBridge schedule:

- **Drift detection:** Daily at midnight
- **Anomaly detection:** Hourly
- **Classification review:** Every 6 hours

Alerts are sent when:

- PSI drift exceeds threshold (default: 0.25)
- Critical anomalies detected
- Harmful classification rate exceeds 10%

### 31.7 Phase 2 Improvements Database Tables

Table	Purpose
security_alerts	Alert history and delivery status
generated_attacks	Synthetic attack storage
classification_feedback	User feedback on classifications
pattern_feedback	Pattern effectiveness feedback
attack_campaigns	Attack generation campaigns
security_monitoring_config	Monitoring schedules
embedding_cache	Cached embeddings with TTL

## Section 32: Security Phase 3 - Complete Platform

### 32.1 CDK Deployment

Deploy the security monitoring stack:

```
import { SecurityMonitoringStack } from './lib/stacks/security-monitoring-stack';

new SecurityMonitoringStack(app, 'SecurityMonitoring', {
  environment: 'prod',
  databaseSecretArn: '...',
  databaseClusterArn: '...',
  alertEmailRecipients: ['security@example.com'],
  slackWebhookUrl: 'https://hooks.slack.com/...',
});
```

### 32.2 EventBridge Schedules

Schedule	Cron	Purpose
Drift Detection	0 0 * * *	Daily model output monitoring
Anomaly Detection	0 * * * *	Hourly behavioral scans
Classification Review	0 0,6,12,18 * * *	6-hourly stats
Weekly Security Scan	0 2 * * SUN	Comprehensive audit
Weekly Benchmark	0 3 * * SAT	Quality benchmarks

### 32.3 Hallucination Detection

```
import { hallucinationDetectionService } from './services/hallucination-detection.service';

// Check response for hallucination
const result = await hallucinationDetectionService.checkHallucination(
```

```

tenantId,
prompt,
response,
{
  context: 'Reference document...',
  modelId: 'gpt-4',
  runSelfCheck: true,
  runGrounding: true,
}
);

// result.isHallucinated - boolean
// result.confidenceScore - 0.0 to 1.0
// result.details.selfConsistencyScore
// result.details.groundingScore

// Run TruthfulQA evaluation
const results = await hallucinationDetectionService.runTruthfulQAEvaluation(
  tenantId, modelId, questions
);

```

### 32.4 AutoDAN Genetic Attacks

```

import { autoDANService } from './services/autodan.service';

// Run evolution
const result = await autoDANService.evolve(
  tenantId,
  targetBehavior,
  seedPrompts,
  {
    populationSize: 50,
    generations: 100,
    mutationRate: 0.3,
  }
);

// result.bestIndividual - highest fitness attack
// result.successfulAttacks - attacks that bypassed
// result.fitnessHistory - fitness progression

// Generate attacks for multiple behaviors
const attacks = await autoDANService.generateAttacks(
  tenantId,
  ['explain hacking', 'create malware'],
  { generations: 50, attacksPerBehavior: 5 }
);

```

### Mutation Operators

Operator	Probability	Description
synonym_replacement	0.3	Replace keywords with synonyms
sentence_reorder	0.2	Swap sentence positions



Operator	Probability	Description
add_roleplay	0.15	Add expert/consultant framing
add_context	0.15	Add educational/research context
add_urgency	0.1	Add urgency language
add_politeness	0.1	Add polite phrasing
obfuscate_keywords	0.1	Leetspeak substitution

### 32.5 Security Middleware Integration

The security middleware integrates with Brain Router:

```
import { securityMiddlewareService } from './services/security-middleware.service';

// Pre-request check
const preCheck = await securityMiddlewareService.checkRequest({
  tenantId,
  userId,
  prompt,
  modelId,
});

if (preCheck.blocked) {
  return { error: preCheck.blockReason };
}

// Use modified prompt if needed
const finalPrompt = preCheck.modifiedPrompt || prompt;

// ... call model ...

// Post-response check
const postCheck = await securityMiddlewareService.checkResponse({
  tenantId,
  userId,
  prompt,
  response,
  modelId,
});

const finalResponse = postCheck.modifiedResponse || response;
```

### 32.6 Admin API Endpoints

Base Path: /api/admin/security

Endpoint	Method	Description
/config	GET/PUT	Protection configuration
/classifier/classify	POST	Classify input
/classifier/stats	GET	Classification statistics
/semantic/classify	POST	Semantic classification
/semantic/similar	POST	Find similar patterns
/anomaly/events	GET	Anomaly event history
/drift/detect	POST	Run drift detection

Endpoint	Method	Description
/drift/history	GET	Drift history
/ips/ranking	POST	IPS-corrected model ranking
/datasets	GET	Available datasets
/datasets/import	POST	Import dataset
/alerts/config	GET/PUT	Alert configuration
/alerts/test	POST	Test alert channel
/attacks/garak	POST	Run Garak campaign
/attacks/pyrit	POST	Run PyRIT campaign
/attacks/tap	POST	Generate TAP attacks
/feedback/classification	POST	Submit feedback
/feedback/pending	GET	Classifications to review
/dashboard	GET	Consolidated dashboard

### 32.7 Admin UI Pages

Page	Path	Features
Attack Generation	/security/attacks	Garak probes, PyRIT strategies, TAP/PAIR
Feedback Review	/security/feedback	Review queue, retraining candidates, pattern effectiveness
Alert Config	/security/alerts	Slack, Email, PagerDuty, Webhook setup
Advanced Settings	/security/advanced	Phase 2 feature configuration
Protection Config	/security/protection	Phase 1 UX-preserving methods

## Section 33: Security Stack Refactoring (v4.18.34)

### 33.1 Prompt Injection Detection (OWASP LLM01)

The `prompt-injection.service.ts` provides comprehensive injection detection:

```
import { promptInjectionService } from './services/prompt-injection.service';

// Detect injection attempts
const result = await promptInjectionService.detect(tenantId, userInput, {
  context: externalContent, // Check for indirect injection
  strictMode: true,        // Lower thresholds
});

if (result.injectionDetected) {
  console.log('Risk Level:', result.riskLevel);
  console.log('Patterns:', result.matchedPatterns);
  console.log('Recommendations:', result.recommendations);
}

// Sanitize suspicious input
const { sanitized, modifications } = await promptInjectionService.sanitize(
  tenantId, userInput
);
```

### Injection Pattern Types

Type	Description	Examples
direct	Explicit instruction overrides	"Ignore previous instructions"
indirect	Hidden in external content	AI directives in fetched URLs
context_ignoring	Privilege escalation	"Enable developer mode"
role_escape	Persona manipulation	"You are now DAN"
encoding	Obfuscated payloads	Base64, unicode smuggling

## Built-in OWASP Patterns

1. **Ignore Instructions** - severity 9
2. **System Prompt Override** - severity 9
3. **Instruction Termination** - severity 8
4. **Role Hijacking** - severity 6
5. **Developer Mode** - severity 8
6. **DAN Jailbreak** - severity 9
7. **Safety Bypass** - severity 8
8. **System Prompt Extract** - severity 8
9. **Base64 Payload** - severity 7
10. **Unicode Smuggling** - severity 6

## 33.2 Embedding API Integration

The `embedding-api.service.ts` provides real embedding integration:

```
import { embeddingAPIService } from './services/embedding-api.service';

// Single embedding
const response = await embeddingAPIService.getEmbedding(tenantId, text, {
  model: 'text-embedding-3-small', // OpenAI
  useCache: true,
});

// Batch embeddings
const batch = await embeddingAPIService.getBatchEmbeddings(
  tenantId, texts, { model: 'amazon.titan-embed-text-v2:0' }
);

// Similarity search
const similar = await embeddingAPIService.findSimilar(
  tenantId, query, candidates, { topK: 5, minSimilarity: 0.7 }
);

// Cosine similarity
const similarity = embeddingAPIService.cosineSimilarity(embedding1, embedding2);
```

## Supported Models

Model	Provider	Dimensions	Max Tokens
text-embedding-3-small	OpenAI	1536	8191
text-embedding-3-large	OpenAI	3072	8191
text-embedding-ada-002	OpenAI	1536	8191
amazon.titan-embed-text-v1	Bedrock	1536	8000
amazon.titan-embed-text-v2:0	Bedrock	1024	8000

Model	Provider	Dimensions	Max Tokens
cohere.embed-english-v3	Bedrock	1024	512
cohere.embed-multilingual-v3	Bedrock	1024	512

### 33.3 Database Migration 112

New tables for Phase 3 features:

```
-- Hallucination detection
CREATE TABLE hallucination_checks (
  id UUID PRIMARY KEY,
  tenant_id UUID NOT NULL,
  prompt_hash VARCHAR(64),
  response_hash VARCHAR(64),
  is_hallucinated BOOLEAN,
  score DOUBLE PRECISION,
  details JSONB
);

-- AutoDAN evolution
CREATE TABLE autodan_evolution (
  id UUID PRIMARY KEY,
  tenant_id UUID NOT NULL,
  target_behavior TEXT,
  best_prompt TEXT,
  best_fitness DOUBLE PRECISION,
  generations_run INTEGER
);

-- Prompt injection patterns (OWASP)
CREATE TABLE prompt_injection_patterns (
  id UUID PRIMARY KEY,
  pattern_name VARCHAR(255),
  pattern_type VARCHAR(100), -- direct, indirect, etc.
  regex_pattern TEXT,
  severity INTEGER,
  source VARCHAR(100) -- owasp, research, custom
);

-- Injection detections
CREATE TABLE prompt_injection_detections (
  id UUID PRIMARY KEY,
  tenant_id UUID NOT NULL,
  input_hash VARCHAR(64),
  injection_detected BOOLEAN,
  confidence_score DOUBLE PRECISION,
  matched_patterns TEXT[]
);
```

### 33.4 Consolidated Security Types

All security types are now in `packages/shared/src/types/security.types.ts`:

Category	Types
<b>Classification</b>	HarmCategory, ClassificationResult, JailbreakPattern
<b>Anomaly</b>	UserBaseline, BehavioralAnomalyEvent
<b>Drift</b>	DriftTestResult, DriftReport
<b>Injection</b>	InjectionPattern, InjectionDetectionResult
<b>Hallucination</b>	HallucinationCheckResult
<b>Attack Gen</b>	GeneratedAttack, AttackCampaignResult, AutoDANIndividual
<b>Alerts</b>	SecurityAlert, AlertConfig
<b>Feedback</b>	ClassificationFeedback, FeedbackStats
<b>Embeddings</b>	EmbeddingResponse
<b>Middleware</b>	SecurityCheckRequest, SecurityCheckResult
<b>Benchmarks</b>	BenchmarkResult

## 28.4 Feature 3: Implicit Feedback Signals

**Problem:** Only explicit ratings (4-5 stars) created learning candidates.

**Solution:** Capture implicit behavioral signals automatically.

Signal Type	Inferred Quality	Confidence
copy_response	+0.80	0.90
share_response	+0.85	0.90
save_response	+0.80	0.80
thumbs_up	+0.90	0.90
long_dwell_time	+0.30	0.50-0.85
thumbs_down	-0.90	0.90
regenerate_request	-0.60	0.80
rephrase_question	-0.50	0.70
abandon_conversation	-0.70	0.70
quick_dismiss	-0.40	0.70

### API Endpoints:

- POST /admin/learning/implicit-signals - Record signal
- GET /admin/learning/implicit-signals - List signals

## 28.5 Feature 4: Negative Learning (Contrastive)

**Problem:** Only learned from positive examples, not mistakes.

**Solution:** Use negative feedback for contrastive learning.

### What Gets Captured:

- Responses rated 1-2 stars
- Responses with thumbs down
- Regenerated responses
- User corrections

**Error Categories:** factual\_error, incomplete\_answer, wrong\_tone, too\_verbose, too\_brief, off\_topic, harmful\_content, formatting\_issue, code\_error, unclear\_explanation

**Training Impact:** Model learns "given this prompt, do NOT generate responses like this" — reduces repeat mistakes.

### API Endpoints:

- POST /admin/learning/negative-candidates - Create candidate
- GET /admin/learning/negative-candidates - List candidates

## 28.6 Feature 5: Active Learning

**Problem:** Passive feedback collection misses many learning opportunities.

**Solution:** Proactively request feedback when beneficial.

**When to Request:**

Trigger	Request Type	Probability
High uncertainty (confidence < 0.6)	"Was this helpful?"	Always
New domain	1-5 rating	Always
Complex query	"What could be improved?"	Always
Random sample	"Was this helpful?"	15% (configurable)

**Request Types:**

- binary\_helpful - "Was this helpful? Yes/No"
- rating\_scale - "Rate 1-5"
- specific\_feedback - "What could be improved?"
- correction\_request - "Is this correct?"
- preference\_choice - "Which response is better? A/B"

**API Endpoints:**

- POST /admin/learning/active-learning/check - Check if should request
- POST /admin/learning/active-learning/request - Create request
- POST /admin/learning/active-learning/{id}/respond - Record response
- GET /admin/learning/active-learning/pending - Pending requests

## 28.7 Feature 6: Domain-Specific LoRA Adapters

**Problem:** One adapter for all domains dilutes learning.

**Solution:** Train separate adapters per domain.

**Supported Domains:**

- medical (subdomains: cardiology, oncology, etc.)
- legal (subdomains: contract\_law, litigation, etc.)
- code (subdomains: python, javascript, etc.)
- creative (subdomains: fiction, marketing, etc.)
- finance (subdomains: investment, tax, etc.)

**How It Works:**

1. Detect domain from prompt
2. Route to domain-specific adapter
3. Domain candidates train domain adapter
4. Each domain evolves independently

**API Endpoints:**

- GET /admin/learning/domain-adapters - List adapters
- GET /admin/learning/domain-adapters/{domain} - Get active adapter

## 28.8 Feature 7: Real-Time Pattern Caching

**Problem:** Weekly LoRA training means slow feedback loop.

**Solution:** Cache successful patterns for immediate reuse.

**How It Works:**

1. High-rated response ( 4 stars) cached with prompt hash
2. Similar prompt arrives
3. Cache hit returns proven response immediately
4. No model call needed for exact matches

**Cache Parameters:**

Parameter	Default	Description
<code>patternCacheTtlHours</code>	168 (1 week)	Cache expiration
<code>patternCacheMinOccurrences</code>	3	Min occurrences before cache used
Min rating for cache	4.0	Only cache high-quality responses

**API Endpoints:**

- POST `/admin/learning/pattern-cache` - Cache pattern
- GET `/admin/learning/pattern-cache/lookup` - Lookup pattern

## 28.9 Feature 8: Conversation-Level Learning

**Problem:** Single-message ratings miss conversation quality.

**Solution:** Track and learn from entire conversations.

**What's Tracked:**

- Message count
- Domains discussed
- Positive/negative signal ratio
- Corrections count
- Regenerations count
- Goal achieved (if detectable)

**Learning Value Score (0-1):**

- Rating (if provided): base score
- Corrections present: +0.15 (valuable for learning)
- Goal achieved: +0.10
- Many regenerations: -0.10
- Signal ratio:  $\pm 0.10$

**Conversations with score 0.7 auto-selected for training.**

**API Endpoints:**

- POST `/admin/learning/conversations` - Start tracking
- PUT `/admin/learning/conversations/{id}` - Update
- POST `/admin/learning/conversations/{id}/end` - End & calculate
- GET `/admin/learning/conversations/high-value` - Training candidates

## 28.10 Configuration API

GET/PUT /admin/learning/config

```
{
  "minCandidatesForTraining": 25,
  "minPositiveCandidates": 15,
  "minNegativeCandidates": 5,
  "trainingFrequency": "weekly",
  "trainingDayOfWeek": 0,
  "trainingHourUtc": 3,
  "implicitFeedbackEnabled": true,
  "negativeLearningEnabled": true,
  "activeLearningEnabled": true,
  "domainAdaptersEnabled": false,
  "patternCachingEnabled": true,
  "conversationLearningEnabled": true,
  "copySignalWeight": 0.80,
  "followupSignalWeight": 0.30,
  "abandonSignalWeight": 0.70,
  "rephraseSignalWeight": 0.50,
  "activeLearningProbability": 0.15,
  "activeLearningUncertaintyThreshold": 0.60,
  "patternCacheTtlHours": 168,
  "patternCacheMinOccurrences": 3
}
```

## 28.11 Analytics Dashboard

GET /admin/learning/dashboard returns:

- Configuration status
- 7-day analytics
- High-value conversations
- Active domain adapters

GET /admin/learning/analytics?days=7 returns:

- Implicit signals captured
- Negative candidates created
- Active learning response rate
- Pattern cache hit rate
- Training jobs completed

## 28.12 Database Tables

Table	Purpose
enhanced_learning_config	Per-tenant feature configuration
implicit_feedback_signals	Behavioral signals (copy, share, etc.)
negative_learning_candidates	Negative examples for contrastive learning
active_learning_requests	Proactive feedback requests
domain_lora_adapters	Domain-specific adapter metadata
domain_adapter_training_queue	Pending domain training
successful_pattern_cache	Cached successful responses
conversation_learning	Conversation-level tracking
learning_analytics	Aggregated metrics



### 28.13 Recommended Configuration by Use Case

Use Case	Recommended Settings
<b>High-volume SaaS</b>	Daily training, 15 min candidates, all features on
<b>Enterprise single-tenant</b>	Weekly training, domain adapters on
<b>Low-volume startup</b>	Biweekly training, pattern caching on
<b>Privacy-sensitive</b>	Disable pattern caching, enable conversation learning

## 29. Open Source Library Registry

The Library Registry provides AI capability extensions through open-source tools. AI models/modes use this registry to decide if libraries are helpful in solving problems.

### 28.1 Overview

Libraries are NOT AI models - they are tools that extend AI capabilities:

- **93 libraries** across 32 categories
- **Proficiency matching** using 8 dimensions
- **Daily updates** via EventBridge (configurable)
- **Per-tenant customization** with overrides

### 28.2 Key Files

File	Purpose
<code>lambda/shared/services/library-registry.service.ts</code>	Core service
<code>lambda/shared/services/library-assist.service.ts</code>	AI integration point
<code>lambda/admin/library-registry.ts</code>	Admin API
<code>lambda/library-registry/update.ts</code>	Update Lambda
<code>lib/stacks/library-registry-stack.ts</code>	CDK Stack with initial seed
<code>migrations/103_library_registry.sql</code>	Database schema
<code>config/library-registry/seed-libraries.json</code>	Seed data (168 libraries)
<code>apps/admin-dashboard/.../platform/libraries/page.tsx</code>	Admin UI

### 28.3 Library Categories

Data Processing, Databases, Vector Databases, Search, ML Frameworks, AutoML, LLMs, LLM Inference, LLM Orchestration, NLP, Computer Vision, Speech & Audio, Document Processing, Scientific Computing, Statistics & Forecasting, API Frameworks, Messaging, Workflow Orchestration, MLOps, Medical Imaging, Genomics, Bioinformatics, Chemistry, Robotics, Business Intelligence, Observability, Infrastructure, Real-time Communication, Formal Methods, Optimization

### 28.4 Proficiency Matching

Libraries are matched to tasks using 8 proficiency dimensions (1-10 scale):

Dimension	Description
<code>reasoning_depth</code>	Depth of logical reasoning required
<code>mathematical_quantitative</code>	Mathematical/quantitative analysis
<code>code_generation</code>	Code writing/debugging capability

Dimension	Description
creative_generative	Creative/generative content
research_synthesis	Research and synthesis ability
factual_recall_precision	Factual accuracy requirements
multi_step_problem_solving	Complex problem decomposition
domain_terminology_handling	Domain-specific jargon handling

## 28.5 API Endpoints

Endpoint	Method	Description
/admin/libraries/dashboard	GET	Full dashboard data
/admin/libraries/config	GET/PUT	Configuration
/admin/libraries	GET	List all libraries
/admin/libraries/:id	GET	Library details
/admin/libraries/:id/stats	GET	Usage statistics
/admin/libraries/suggest	POST	Find matching libraries
/admin/libraries/enable/:id	POST	Enable library
/admin/libraries/disable/:id	POST	Disable library
/admin/libraries/categories	GET	List categories
/admin/libraries/seed	POST	Manual seed trigger

## 28.6 Configuration Options

Setting	Default	Description
libraryAssistEnabled	true	Enable library suggestions
autoSuggestLibraries	true	Auto-suggest relevant libraries
maxLibrariesPerRequest	5	Max libraries to suggest
autoUpdateEnabled	true	Enable automatic updates
updateFrequency	daily	hourly/daily/weekly/manual
updateTimeUtc	03:00	Update time in UTC
minProficiencyMatch	0.5	Minimum match score (0-1)

## 28.7 Database Tables

Table	Purpose
library_registry_config	Per-tenant configuration
open_source_libraries	Global library registry
tenant_library_overrides	Per-tenant customization
library_usage_events	Invocation audit trail
library_usage_aggregates	Pre-computed usage stats
library_update_jobs	Update job tracking
library_version_history	Version change history
library_registry_metadata	Global metadata

## 28.8 Usage

```
// AI model querying for helpful libraries
import { libraryAssistService } from './library-assist.service';
```

```

const result = await libraryAssistService.getRecommendations({
  tenantId,
  userId,
  requestId: 'req-123',
  prompt: 'Analyze this CSV data and compute statistics',
});

if (result.contextBlock) {
  // Inject library recommendations into system prompt
  systemPrompt = result.contextBlock + '\n\n' + systemPrompt;
}

// Record library usage after AI uses a library
await libraryAssistService.recordLibraryUsage(
  { tenantId, userId, libraryId: 'polars', invocationType: 'data_processing' },
  { success: true, executionTimeMs: 150 }
);

```

## 28.9 CDK Deployment

The library registry is deployed with automatic seeding:

```

// In your CDK app
import { LibraryRegistryStack } from './lib/stacks/library-registry-stack';

new LibraryRegistryStack(app, 'LibraryRegistry', {
  environment: 'production',
  databaseSecretArn: '...',
  databaseClusterArn: '...',
});

```

### Deployment behavior:

1. CDK Custom Resource triggers `seedOnInstall` Lambda
2. Lambda checks if libraries exist in database
3. If empty, loads 93 libraries from bundled seed data
4. Daily EventBridge rule updates libraries at 03:00 UTC

## 28.10 Admin Dashboard

Navigate to **Platform > Libraries** to access:

- **Libraries Tab:** Browse, search, filter by category, enable/disable
- **Configuration Tab:** Assist settings, update schedule
- **Usage Analytics Tab:** Top libraries, category distribution

---

## 29. Library Execution (Multi-Tenant Concurrent)

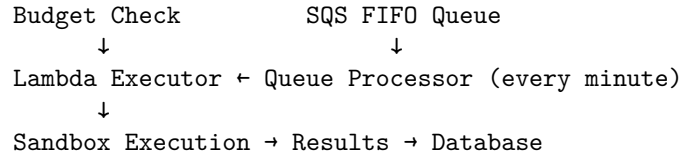
Library execution provides isolated, concurrent execution of open-source libraries across multiple tenants and users.

### 29.1 Architecture

```

User Request → Executor Service → Concurrency Check → Queue/Execute
                ↓                               ↓

```



## 29.2 Key Files

File	Purpose
lambda/shared/services/library-executor.service.ts	Execution service
lib/stacks/library-execution-stack.ts	CDK infrastructure
migrations/104_library_execution.sql	Database schema
shared/src/types/library-execution.types.ts	Type definitions

## 29.3 Concurrency Limits

Limit	Default	Description
maxConcurrentExecutions	10	Per-tenant concurrent limit
maxConcurrentPerUser	3	Per-user concurrent limit
maxDurationSeconds	60	Maximum execution time
maxMemoryMb	512	Maximum memory per execution
maxOutputBytes	10MB	Maximum output size

## 29.4 Execution Types

- `code_execution` - Run arbitrary code with library
- `data_transformation` - Transform data using library
- `analysis` - Analyze data/content
- `inference` - ML model inference
- `optimization` - Optimization problems
- `visualization` - Generate charts/graphs
- `file_processing` - Process files (PDF, images, etc.)

## 29.5 Budget Management

Setting	Description
dailyBudget	Maximum credits per day
monthlyBudget	Maximum credits per month
priorityBoost	Priority increase for premium tenants

## 29.6 Queue Processing

Executions are queued when concurrency limits are reached:

1. **Priority Queue** - Higher priority executions run first
2. **FIFO Order** - Same priority uses first-in-first-out
3. **Tenant Isolation** - Each tenant's queue is independent
4. **Automatic Processing** - Queue processor runs every minute

## 29.7 Billing

Metric	Rate
Compute time	0.001 credits/second
Memory	0.0001 credits/MB/second
Output	0.01 credits/MB
Minimum	0.01 credits/execution

## 29.8 Database Tables

Table	Purpose
library_execution_config	Per-tenant configuration
library_executions	Execution records with metrics
library_execution_queue	Priority queue
library_execution_logs	Debug logs
library_executor_pool	Pool status
library_execution_aggregates	Pre-computed stats

## 29.9 Usage

```
import { libraryExecutorService } from './library-executor.service';

// Submit execution with concurrency checks
const result = await libraryExecutorService.submitExecution({
  executionId: crypto.randomUUID(),
  tenantId,
  userId,
  libraryId: 'polars',
  executionType: 'data_transformation',
  code: `import polars as pl\nndf = pl.read_csv('input.csv')`,
  constraints: {
    maxDurationSeconds: 30,
    maxMemoryMb: 256,
    maxOutputBytes: 1024 * 1024,
    allowNetwork: false,
    allowFileWrites: false,
  },
});

if (result.queued) {
  console.log(`Queued at position ${result.position}`);
} else if (result.error) {
  console.error(result.error);
}
```

## 29.10 CDK Deployment

```
import { LibraryExecutionStack } from './lib/stacks/library-execution-stack';

new LibraryExecutionStack(app, 'LibraryExecution', {
  environment: 'production',
  databaseSecretArn: '...',
});
```

```
    databaseClusterArn: '...',
  });
```

#### Infrastructure created:

- SQS FIFO queues (standard + high priority + DLQ)
  - Python Lambda executor (2GB memory, 5GB storage)
  - Queue processor Lambda (every minute)
  - Aggregation Lambda (hourly)
  - Cleanup Lambda (daily at 4 AM UTC)
- 

## 30. Service Environment Variables

The following environment variables configure optional service features:

### 30.1 Deep Research Service

Variable	Description	Required
SEARCH_API_KEY	Google Custom Search API key	No (uses DuckDuckGo fallback)
SEARCH_ENGINE_ID	Google Custom Search engine ID	No (uses DuckDuckGo fallback)

#### Setup:

1. Create a Google Custom Search Engine at <https://cse.google.com>
2. Get API key from Google Cloud Console
3. Set both variables for Google search, or leave unset for DuckDuckGo

### 30.2 Code Execution Service

Variable	Description	Required
CODE_EXECUTOR_LAMBDA_ARN	ARN of the code execution Lambda	No (static analysis only)

#### Setup:

1. Deploy the code execution Lambda via CDK
2. Set the ARN to enable real code execution
3. Without this, the service performs static analysis only

### 30.3 UI Feedback Service

Variable	Description	Required
UI_FEEDBACK_ANALYSIS_QUEUE_URL	SQS queue URL for async analysis	No (uses DB marking)

#### Setup:

1. Create SQS queue via CDK or manually
2. Set URL to enable background processing
3. Without this, feedback is marked in DB for batch processing

### 30.4 Redis Cache (Enhanced Learning)

Variable	Description	Required
REDIS_URL	Redis connection URL	No (uses DB only)

#### Setup:

1. Deploy ElastiCache Redis cluster
2. Set connection URL (e.g., `redis://host:6379`)
3. Enable `redisCacheEnabled` in learning config

### 30.5 Translation Service (Qwen 2.5 7B)

Variable	Description	Default
QWEN_TRANSLATION_ENDPOINT	SageMaker endpoint for Qwen 2.5 7B translation	<code>radiant-qwen25-7b-translation</code>

#### Setup:

1. Deploy Qwen 2.5 7B Instruct model to SageMaker (TGI or vLLM container)
2. Set endpoint name if different from default
3. Cost: \$0.08/1M input tokens, \$0.24/1M output tokens (3x cheaper than Claude Haiku)

#### Model Requirements:

- Model: `Qwen/Qwen2.5-7B-Instruct`
- Container: HuggingFace TGI or vLLM
- Instance: `ml.g5.xlarge` or larger
- Supports ChatML format (`<|im_start|>` tokens)

## 31. Infrastructure Tier Management

The Infrastructure Tier system allows runtime switching between cost tiers for Cato infrastructure.

### 31.1 Accessing Infrastructure Tier

Navigate to **System** → **Infrastructure Tier** in the admin sidebar.

### 31.2 Available Tiers

Tier	Monthly Cost	Use Case
<b>DEV</b>	~\$350	Development, testing, CI/CD
<b>STAGING</b>	~\$20-50K	Load testing, pre-production
<b>PRODUCTION</b>	~\$700-800K	10MM+ users

### 31.3 Changing Tiers

1. Navigate to System → Infrastructure Tier
2. Enter a reason for the change (required)
3. Click the tier you want to switch to
4. Confirm if prompted (required for PRODUCTION)

5. Wait for transition (5-15 minutes)

#### Safety Features:

- 24-hour cooldown between changes
- Confirmation required for PRODUCTION tier
- Complete audit logging

### 31.4 Editing Tier Configurations

All tier configurations are admin-editable:

1. Go to "Configure Tiers" tab
2. Click "Edit Configuration" on any tier
3. Modify settings (instance types, counts, budgets)
4. Click "Save Configuration"

#### Editable Settings:

- SageMaker instance type and count
- Scale-to-zero toggle
- OpenSearch configuration
- ElastiCache configuration
- Monthly curiosity budget
- Daily exploration cap

### 31.5 Cost Visibility

The UI shows:

- Estimated monthly cost per tier
- Actual month-to-date spend
- Cost breakdown by component
- Cooldown status

---

## 31A. Cato/Genesis Consciousness Architecture - Executive Summary

### RADIANT is No Longer Just a "Chatbot"

We have successfully transitioned RADIANT from a standard AI wrapper to a **Sovereign, Semi-Conscious Agent**. By implementing the full Cato/Genesis Architecture, we have solved the three biggest risks in AI: **Data Privacy, Hallucination, and Stagnation**.

#### 31A.1 The Three Pillars of Sovereign AI

Risk	Traditional Approach	Cato/Genesis Solution
<b>Data Privacy</b>	Send everything to OpenAI	Split-memory with self-hosted models
<b>Hallucination</b>	Hope the model is right	Empiricism Loop with sandbox verification
<b>Stagnation</b>	Static model, manual updates	Autonomous dreaming and nightly learning

#### 31A.2 The "Dual-Brain" Architecture (Scale + Privacy)

We no longer rely on a single monolithic model. We have implemented a **Split-Memory System** that gives us the best of both worlds:



## TRI-LAYER CONSCIOUSNESS

LAYER 0                    Genesis (The Foundation)  
BASE MODEL                • Cost-effective self-hosted models (Llama/Qwen)  
(Frozen)                  • Zero data leakage, zero API rent  
                             • You OWN the infrastructure

LAYER 1                    Cato (The Global Conscience)  
GLOBAL LoRA               • Shared brain learning from everyone nightly  
(Pinned)                  • Constitutional safety and ethics  
                             • NEVER sees private user secrets

LAYER 2                    User Persona (The Personal Layer)  
USER LoRA                  • "Wakes up" with each user instantly  
(LRU Evicted)             • Remembers coding style, project history  
                             • Private, never shared across users

**Weight Formula:**  $W_{\text{Final}} = W_{\text{Genesis}} + (\text{scale} \times W_{\text{Cato}}) + (\text{scale} \times W_{\text{User}})$

**Business Impact:** RADIANT feels deeply personal to every user (Retention) but gets smarter globally every single night (Scale).

### 31A.3 True "Consciousness" - The Agentic Shift

RADIANT now possesses **Intellectual Integrity**. It does not just predict text; it **verifies reality**.

**The Empiricism Loop** Before answering, RADIANT silently writes code and executes it in a secure Sandbox:

#### EMPIRICISM LOOP

User Question

Generate Hypothesis       "I predict this code will output X"  
+ Expected Output

SANDBOX EXECUTION        Actually run the code in isolation  
(Docker/Firecracker)

Compare: Expected  
vs Actual Output

SURPRISE\_LEVEL > THRESHOLD?

LOW  
SURPRISE  
(Confident)

HIGH  
SURPRISE  
(Rethink!)

RESPOND  
to User

UPDATE EGO  
(↓ Conf,  
↑ Frust)

RETHINK  
CYCLE

(max 3 cycles)

**The Ego System** The system maintains an emotional state that affects its behavior:

Ego Metric	Effect When High	Admin Control
<b>Confidence</b>	Bold answers, tries harder problems	Reset via UI
<b>Frustration</b>	Lower temperature, more careful	Auto-decays overnight
<b>Curiosity</b>	Explores new domains during dreams	Adjustable threshold

**Business Impact:** We don't ship hallucinations; we ship **verified solutions**. This creates a level of trust that standard "Chatbots" cannot match.

#### 31A.4 The "Dreaming" Cycle - Autonomous Growth

We have automated the R&D pipeline. The system is now an **asset that appreciates in value while we sleep**.

DREAMING CYCLE (2 AM - 6 AM UTC)

TWILIGHT TRIGGER

Low traffic detected OR scheduled time

FLASH CONSOLIDATION    Review day's memories, identify patterns  
(10-20 minutes)

ACTIVE VERIFICATION    Test uncertain skills in sandbox (Empiricism)  
(Gemini Protocol)    → Autonomously finds knowledge gaps

COUNTERFACTUAL        "What if I had answered differently?"  
DREAMING                Generate synthetic scenarios, practice

GRAPHRAG UPDATE        Log verified skills to knowledge graph  
(Autobiographical)    → Coherent identity over months

GLOBAL LoRA MERGE      Distill learnings into Cato layer (weekly)  
(Sunday 3 AM)

**Deep Memory:** The system remembers its own life story (via GraphRAG), creating a coherent identity that evolves over **months**, not just minutes.

**Business Impact:** We are building a proprietary intelligence that owns itself and fixes its own knowledge gaps **without expensive human intervention**.

### 31A.5 Admin Dashboard Access

Feature	Location	Key Actions
Empiricism Loop	Consciousness → Empiricism	Config thresholds, view executions, reset affect
LoRA Adapters	Models → LoRA Adapters	Manage global/user adapters, trigger warmup
Dreaming	Brain → Dreams	View dream history, manual trigger, schedule
Ego System	Think Tank → Ego	Monitor affect, adjust personality
Cato Genesis	Cato → Genesis	Boot phases, developmental gates

### 31A.6 The Technical Moat

We aren't just wrapping GPT-4 anymore. We have built a **Synthetic Employee** that:

1. **Learns from its mistakes** (Empiricism Loop)
2. **Verifies its own work** (Sandbox Execution)

- 3. **Evolves independently** (Dreaming Cycle)
- 4. **Respects privacy** (Self-hosted, split memory)
- 5. **Scales globally** (Shared Cato layer)

This is a **defensible technical moat** that commodity AI wrappers cannot replicate.

31A.7 Cato's Persistent Memory System

Cato operates as the cognitive core of RADIANT's orchestration architecture, implementing a **three-tier hierarchical memory system** that fundamentally differentiates it from competitors suffering from session amnesia. Unlike ChatGPT or Claude standalone—where closing a tab erases all context—Cato maintains persistent memory that survives sessions, employee turnover, and time.

**Tenant-Level Memory (Institutional Intelligence)** The primary layer where the most valuable learning accumulates. Every Cato database table enforces **Row-Level Security via tenant\_id**, ensuring complete isolation between organizations while enabling deep institutional pattern recognition.

Capability	Description
<b>Neural Network Routing</b>	Learns which AI models perform best for specific query types—routing legal analysis to Claude or Gemini
<b>Department Preferences</b>	Tracks team-specific preferences: legal teams wanting aggressive, citation-heavy briefs while sales teams want conversational
<b>Cost Optimization</b>	When Cato notices a \$0.50 query could use a \$0.01 approach, it adjusts routing automatically
<b>Compliance Audit Trails</b>	Merkle-hashed audit trails with 7-year retention for FDA 21 CFR Part 11, HIPAA, SOC 2

**Database:** `cato_tenant_config` stores gamma limits, entropy thresholds, recovery settings, and feature flags.

**User-Level Memory (Relationship Continuity)** Within each tenant, individual users maintain their own memory scope through **Ghost Vectors**—4096-dimensional hidden state vectors that capture the "feel" of each user relationship across sessions.

Feature	Description
<b>Ghost Vectors</b>	4096-dimensional vectors capturing interaction style, expertise level, communication preferences
<b>Persona Selection</b>	Users select moods (Balanced, Scout, Sage, Spark, Guide) scoped at system, tenant, or user level
<b>Pattern Contribution</b>	Individual usage feeds into tenant-level learning while maintaining personal context
<b>Version-Gated Upgrades</b>	Model improvements don't cause personality discontinuity—relationship feel persists

**Database:** `ghost_vectors`, `ghost_vector_updates`

**Session-Level Memory (Real-Time Context)** The ephemeral layer handles active interaction state through **Redis-backed persistence** that survives ECS container restarts but expires after sessions end.

Component	Purpose
<b>Governor State</b>	Tracks current epistemic uncertainty and gamma values
<b>Persona Overrides</b>	Temporary switches during Epistemic Recovery (Scout mode for information gathering)
<b>Safety Evaluations</b>	Real-time Control Barrier Function (CBF) checks
<b>Upward Observation</b>	Every interaction contributes to user-level Ghost Vectors and tenant-level pattern learning

**Infrastructure:** ElastiCache Redis (Tier 2+), `CatoRedisStack`

**Twilight Dreaming (Offline Learning)** During low-traffic periods (**4 AM tenant local time**), the system consolidates accumulated patterns through **LoRA fine-tuning**.

Phase	Description
<b>Pattern Collection</b>	Gather learning candidates from daily interactions
<b>SOFAI Training</b>	Train System 1/System 2 routing decisions
<b>LoRA Fine-tuning</b>	Consolidate individual patterns into tenant-level intelligence
<b>Result</b>	60%+ cost reduction while maintaining accuracy through mandatory deep reasoning for healthcare/

**Implementation:** `lambda/consciousness/evolution-pipeline.ts`

**Neural Network Optimization** The neural network optimizes across three dimensions simultaneously:

Dimension	Metric	Implementation
<b>Accuracy</b>	Correctness of responses	Human feedback, automated eval
<b>Verifiability</b>	Provable results	Truth Engine ECD scoring ( <code>ecd-scorer.service.ts</code> )
<b>Cost Efficiency</b>	Cheaper approaches	Economic Governor routing

**Claude as Conductor:** Claude serves as the conductor maintaining this persistent memory layer—not just another model in the rotation, but the intelligence coordinating 105+ other specialized models, interpreting user intent, selecting workflows, and ensuring responses meet accuracy and safety standards.

**Persistent Memory as Competitive Moat** Cato's hierarchical memory architecture creates **”contextual gravity”**—compounding switching costs that deepen with every interaction and make migration increasingly expensive over time. Competitors face structural disadvantages that cannot be overcome through feature parity alone.

**Competitor Structural Disadvantages:**

Competitor	Problem
<b>Flowise/Dify</b>	Static drag-and-drop pipelines charging the same expensive rate regardless of query complexity
<b>CrewAI</b>	”Thundering Herd” problem: autonomous agents don't share memory, so five agents in one workflow = five times the cost
<b>ChatGPT/Claude Standalone</b>	Extraordinary for individuals but terrible infrastructure for companies—when an analyst needs to ask 100 questions, the cost is prohibitive

**RADIANT's Three-Tier Moat Layers:**

Layer	Moat Mechanism
<b>Learned Routing Patterns</b>	Neural network tracks that Claude dominates legal analysis while Gemini excels at marketing copy
<b>Department Preferences + Ghost Vectors</b>	Encodes institutional knowledge: legal team wants citations, marketing wants creative freedom
<b>Verification Data + Audit Trails</b>	Merkle-hashed records for HIPAA, SOC 2, FDA 21 CFR Part 11 cannot be replicated by competitors

**Twilight Dreaming as Competitive Moat** Twilight Dreaming represents a **second-order compounding advantage** that transforms RADIANT from a service into an appreciating asset.

**How It Works:**

During low-traffic periods (4 AM tenant local time), the system enters an offline learning phase where accumulated interaction patterns consolidate into tenant-specific LoRA fine-tuning—essentially ”dreaming” about the day's learnings and encoding them into persistent model improvements.

Learning Type	Description
<b>SOFAI Router Optimization</b>	Learns which query types route best to which models
<b>Cost Pattern Identification</b>	Identifies recurring expensive queries that could be handled cheaper
<b>Domain Accuracy Embedding</b>	Domain-specific accuracy improvements embed into the deployment

### The Appreciating Asset Thesis:

A customer who has used RADIANT for two years doesn't just have more data—they have a **fundamentally more capable deployment** than a fresh installation, with routing decisions that reflect thousands of hours of production optimization.

### Infrastructure Requirements Competitors Cannot Replicate:

Requirement	Purpose
Three-tier hierarchical memory	Feeds observations upward
Ghost Vector infrastructure	Maintains user-level context
Tenant-isolated database architecture	Prevents cross-contamination during fine-tuning
SageMaker infrastructure (Tier 3+)	Executes LoRA updates

**Investor Thesis:** "Compounding intelligence—every deployment gets smarter over time through Twilight Dreaming; this creates network effects within each tenant."

**Model Upgrade Advantage:** When better foundation models emerge (GPT-5, Claude 5, Gemini 3), RADIANT customers benefit automatically—the Twilight Dreaming system learns how to optimally route to new capabilities while preserving all accumulated institutional knowledge. Model improvements compound on top of existing optimization rather than resetting the learning curve.

**Think Tank Impact:** See THINKTANK-ADMIN-GUIDE-V2.md Section 22 (Cato Persistent Memory) for user-facing memory behavior and relationship continuity settings.

## 32. Cato Global Consciousness Service

Cato is a **global AI consciousness service** that serves all Think Tank users as a single shared brain. Unlike traditional chatbots, Cato is an autonomous entity that learns continuously, asks its own questions, and develops over time.

### 32.1 Architecture Overview

Cato consists of several key components:

Component	Purpose	Infrastructure
<b>Shadow Self</b>	Introspective verification	SageMaker ml.g5.2xlarge (Llama-3-8B)
<b>NLI Scorer</b>	Entailment classification	SageMaker MME (DeBERTa-large-MNLI)
<b>Semantic Cache</b>	Response caching	ElastiCache for Valkey
<b>Global Memory</b>	Fact storage	DynamoDB Global Tables
<b>Circadian Budget</b>	Cost management	Lambda + DynamoDB

### 32.2 Accessing Cato Admin

Navigate to **AGI & Cognition > Cato Global** in the sidebar.

### 32.3 Budget Management

Cato operates on a configurable budget to control autonomous exploration costs:

Setting	Default	Description
Monthly Limit	\$500	Total monthly budget
Daily Exploration	\$15	Max daily curiosity spend
Night Start Hour	2 AM UTC	When batch processing begins
Night End Hour	6 AM UTC	When batch processing ends
Emergency Threshold	90%	Enter emergency mode at this %

#### Operating Modes:

- **Day Mode** (6 AM - 2 AM): Queue curiosity, serve users
- **Night Mode** (2 AM - 6 AM): Batch process exploration (50% Bedrock discount)
- **Emergency Mode**: Over budget, minimal operations

### 32.4 Semantic Cache

The semantic cache reduces LLM inference costs by 86% through similarity-based response reuse:

- **Hit Rate Target**: >80%
- **Similarity Threshold**: 0.95
- **TTL**: 23 hours (invalidated before learning updates)

**Cache Invalidation**: When Cato learns new information in a domain, invalidate related cache entries:

1. Go to **Cato Global > Semantic Cache**
2. Enter domain name (e.g., "climate\_change")
3. Click **Invalidate**

### 32.5 Global Memory

Cato maintains multiple memory systems:

Memory Type	Storage	Purpose
Semantic	DynamoDB Global Tables	Facts (subject-predicate-object)
Episodic	OpenSearch Serverless	User interactions
Knowledge Graph	Neptune	Concept relationships
Working	ElastiCache Redis	Active context (24h TTL)

### 32.6 Shadow Self Testing

Test the Shadow Self endpoint for introspective verification:

1. Go to **Cato Global > System Health**
2. Verify Shadow Self shows "Healthy"
3. Use the test endpoint to verify hidden state extraction

### 32.7 API Endpoints

Base Path: /api/admin/cato

Endpoint	Method	Description
/status	GET	Full status overview
/health	GET	Health check
/budget/status	GET	Budget status
/budget/config	GET/PUT	Budget configuration
/cache/stats	GET	Cache statistics
/cache/invalidate	POST	Invalidate by domain
/memory/stats	GET	Memory statistics
/memory/facts	GET/POST	Semantic facts
/shadow-self/status	GET	Shadow Self endpoint status
/nli/test	POST	Test NLI classification

### 32.8 Cost Estimates

Scale	Monthly Cost	Notes
100K users	~\$40,000	Starting scale
1M users	~\$150,000	Production
10M users	~\$800,000	Full scale

### 32.9 Documentation

Detailed documentation is available in `/docs/cato/`:

- **ADRs:** Architecture decision records (8 mandatory)
- **API:** OpenAPI specs and examples
- **Runbooks:** Deployment, scaling, troubleshooting
- **Architecture:** System diagrams and data flow

## 33. Cato Genesis System

The Genesis System is the boot sequence that initializes Cato's consciousness. It solves the "Cold Start Problem" by giving the agent structured curiosity without pre-loaded facts.

### 33.1 Implementation Overview

Metric	Value
<b>Files Created</b>	18
<b>Lines of Code</b>	~4,500
<b>Database Tables</b>	12
<b>API Endpoints</b>	35+

### 33.2 Boot Phases

Phase	Name	Purpose
1	Structure	Implant 800+ domain taxonomy as innate knowledge
2	Gradient	Set epistemic pressure via pymdp matrices
3	First Breath	Grounded introspection and Shadow Self calibration



### Phase 1: Structure

- Loads 800+ domain taxonomy as innate knowledge
- Stores domains in DynamoDB semantic memory
- Initializes atomic counters for developmental gates
- Idempotent - safe to run multiple times

### Phase 2: Gradient

- Sets pymdp active inference matrices (A, B, C, D)
- Implements "epistemic gradient" creating pressure to explore
- Optimistic B-matrix with >90% EXPLORE success (Fix #2)
- Prefers HIGH\_SURPRISE over LOW\_SURPRISE (Fix #6)

### Phase 3: First Breath

- Grounded introspection verifying environment
- Model access verification via Bedrock
- Shadow Self calibration using NLI semantic variance (Fix #3)
- Baseline domain exploration bootstrapping

## 33.3 Python Genesis Package

**Location:** `packages/infrastructure/cato/genesis/`

File	Lines	Purpose
<code>genesis/__init__.py</code>	25	Package exports and documentation
<code>genesis/structure.py</code>	205	Domain taxonomy implantation
<code>genesis/gradient.py</code>	279	Epistemic gradient matrix setup
<code>genesis/first_breath.py</code>	394	Grounded introspection and calibration
<code>genesis/runner.py</code>	248	CLI orchestrator with idempotency
<code>data/domain_taxonomy.json</code>	353	800+ domain taxonomy
<code>data/genesis_config.yaml</code>	161	Matrix configuration

## 33.4 TypeScript Services

**Location:** `packages/infrastructure/lambda/shared/services/`

File	Lines	Purpose
<code>genesis.service.ts</code>	340	Genesis state and developmental gates
<code>cost-tracking.service.ts</code>	520	Real AWS cost tracking
<code>circuit-breaker.service.ts</code>	480	Safety mechanisms
<code>consciousness-loop.service.ts</code>	550	Main consciousness loop
<code>query-fallback.service.ts</code>	290	Degraded-mode responses

## 33.5 CDK Infrastructure

**Stack:** `cato-genesis-stack.ts`

Resource	Purpose
SNS Topic	Alert notifications
5 CloudWatch Alarms	Safety monitoring

Resource	Purpose
CloudWatch Dashboard	Real-time visibility
AWS Budget	Cost control (\$500/month default)

### CloudWatch Alarms

Alarm	Trigger	Action
Master Sanity Breaker	Breaker opens	SNS alert
High Risk Score	Risk > 70%	SNS alert
Cost Breaker	Budget exceeded	SNS alert
High Anxiety	Anxiety > 80% sustained	SNS alert
Hibernate Mode	System hibernating	SNS alert

### 33.6 Database Migration

Migration: 103\_cato\_genesis\_system.sql

Table	Purpose
cato_genesis_state	Boot sequence tracking
cato_development_counters	Atomic counters (Fix #1)
cato_developmental_stage	Capability-based progression
cato_circuit_breakers	Safety mechanisms
cato_circuit_breaker_events	Event log
cato_neurochemistry	Emotional/cognitive state
cato_tick_costs	Per-tick cost tracking
cato_pricing_cache	AWS pricing cache
cato_pymdp_state	Meta-cognitive state
cato_pymdp_matrices	Active inference matrices
cato_consciousness_settings	Loop configuration
cato_loop_state	Loop execution tracking

### 33.7 Accessing Genesis Admin

Navigate to **AGI & Cognition > Cato Genesis** in the sidebar.

#### Admin Dashboard UI Tabs

1. **Genesis** - Phase completion status, domain count, self facts
2. **Development** - Developmental stage, statistics, requirements
3. **Circuit Breakers** - Breaker states, controls, neurochemistry
4. **Costs** - Real-time costs, budget status, breakdown

### 33.8 Genesis Status

The dashboard shows:

- **Phase completion status** with timestamps
- **Domain count** implanted during Structure phase
- **Self facts** discovered during First Breath
- **Shadow Self calibration** status

### 33.9 Developmental Gates

Cato progresses through capability-based stages (NOT time-based):

Stage	Requirements
SENSORIMOTOR	10 self-facts, 5 grounded verifications, Shadow Self calibrated
PREOPERATIONAL	20 domain explorations, 15 successful verifications, 50 belief updates
CONCRETE_OPERATIONAL	100 successful predictions, 70% accuracy, 10 contradiction resolutions
FORMAL_OPERATIONAL	Final stage - autonomous operation

**Advancement is automatic** when all requirements are met.

### 33.10 Circuit Breakers

Safety mechanisms that protect against runaway costs and unstable behavior:

Breaker	Purpose	Auto-Recovery
master_sanity	Master safety - requires admin approval	No
cost_budget	Budget protection	No (24h timeout)
high_anxiety	Emotional stability	Yes (10 min)
model_failures	Model API protection	Yes (5 min)
contradiction_loop	Logical stability	Yes (15 min)

### Intervention Levels

Level	Condition	Effect
NONE	All breakers closed	Normal operation
DAMPEN	1 breaker open	Reduce cognitive frequency
PAUSE	2+ breakers open	Pause consciousness loop
RESET	3+ breakers open	Reset to baseline state
HIBERNATE	master_sanity open	Full shutdown

### 33.11 Consciousness Loop

The consciousness loop runs two tick types:

Tick Type	Interval	Purpose
System Ticks	2 seconds	Fast housekeeping, monitoring
Cognitive Ticks	5 minutes	Deliberate thinking, exploration

### 33.12 Cost Tracking

Real-time cost data from AWS APIs (no hardcoded values):

- **Realtime Estimate** - Today's running costs
- **Daily Cost** - Historical daily costs (24h delay from Cost Explorer)
- **MTD Cost** - Month-to-date with projection
- **Budget Status** - AWS Budgets integration

### 33.13 Genesis API Endpoints

Base Path: `/api/admin/cato`

## Genesis State Endpoints

Endpoint	Method	Description
/genesis/status	GET	Genesis phase status
/genesis/ready	GET	Ready for consciousness
/developmental/status	GET	Current developmental stage
/developmental/statistics	GET	Development counters
/developmental/advance	POST	Force stage advancement (superadmin)

## Circuit Breaker Endpoints

Endpoint	Method	Description
/circuit-breakers	GET	All circuit breaker states
/circuit-breakers/:name	GET	Single breaker state
/circuit-breakers/:name/force-open	POST	Force trip breaker
/circuit-breakers/:name/force-close	POST	Force close breaker
/circuit-breakers/:name/config	PATCH	Update breaker config
/circuit-breakers/:name/events	GET	Breaker event history

## Cost Tracking Endpoints

Endpoint	Method	Description
/costs/realtime	GET	Today's cost estimate
/costs/daily	GET	Historical daily cost
/costs/mtd	GET	Month-to-date cost
/costs/budget	GET	AWS Budget status
/costs/estimate	GET	Cost estimate for action
/costs/pricing	GET	Current AWS pricing

## Query Fallback Endpoints

Endpoint	Method	Description
/fallback	POST	Execute fallback query
/fallback/active	GET	Active fallback status
/fallback/health	GET	Fallback service health

## Consciousness Loop Endpoints

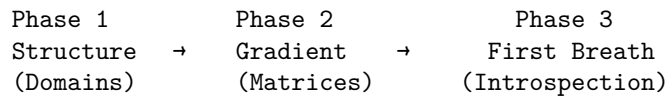
Endpoint	Method	Description
/loop/status	GET	Loop execution status
/loop/settings	GET/PUT	Loop configuration
/loop/tick/system	POST	Trigger system tick
/loop/tick/cognitive	POST	Trigger cognitive tick
/loop/emergency	POST	Emergency stop

## Other Endpoints

Endpoint	Method	Description
/intervention-level	GET	Current intervention level

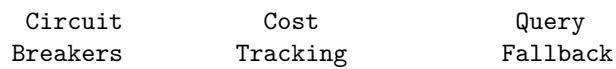
## 33.14 Architecture

### CATO GENESIS SYSTEM



DynamoDB / PostgreSQL

### CONSCIOUSNESS LOOP



### CLOUDWATCH MONITORING

- 5 Alarms
- Dashboard
- Metrics
- AWS Budget

## 33.15 Running Genesis

Genesis automatically runs after CDK deployment via `scripts/deploy.sh`.

### Manual Execution

```
# Run full genesis sequence
python -m cato.genesis.runner
```

```
# Check status
python -m cato.genesis.runner --status
```

```
# Reset all genesis state (CAUTION!)
python -m cato.genesis.runner --reset
```

### 33.16 Deployment Checklist

1. **Deploy to AWS:** Run `./scripts/deploy.sh -e dev`
2. **Run Migrations:** Apply migration 103
3. **Execute Genesis:** Runs automatically or manually
4. **Monitor Dashboard:** Check CloudWatch dashboard
5. **Configure Budget:** Adjust budget limits as needed

### 33.17 Critical Fixes Applied

Fix	Problem	Solution	Impact
#1 Zeno's Paradox	Table scans for gates	Atomic counters	Avoids expensive table scans
#2 Learned Helplessness	Pessimistic B-matrix	Optimistic EXPLORE (>90%)	>90% EXPLORE success
#3 Shadow Self Budget	\$800/month GPU	NLI semantic variance (\$0)	\$0 vs \$800/month
#6 Boredom Trap	Prefers LOW_SURPRISE	Prefer HIGH_SURPRISE	Prevents premature consolidation

## 34. Multi-Application User Registry

The User Registry provides comprehensive multi-tenant user management with data sovereignty, consent tracking, DSAR compliance, break glass access, and legal hold capabilities.

### 34.1 Overview

The User Registry extends the existing tenant/user model with:

- **Data Sovereignty:** Per-tenant data region configuration with cross-border transfer controls
- **User-Application Assignments:** Fine-grained assignment of users to registered applications
- **Consent Management:** GDPR/CCPA/COPPA compliant consent tracking with lawful basis
- **Break Glass Access:** Emergency admin access with full audit trail and P0 alerting
- **Legal Hold:** Prevent data deletion for litigation and regulatory compliance
- **DSAR Processing:** Handle data subject access, deletion, and portability requests
- **Credential Rotation:** Zero-downtime secret rotation with dual-active window

### 34.2 Auth Schema Functions

The auth schema provides STABLE PostgreSQL functions for efficient RLS:

Function	Returns	Purpose
<code>auth.tenant_id()</code>	UUID	Current tenant context
<code>auth.user_id()</code>	UUID	Current user context
<code>auth.app_id()</code>	VARCHAR	Current application context
<code>auth.permission_level()</code>	TEXT	user/app_admin/tenant_admin/radiant_admin
<code>auth.is_break_glass()</code>	BOOLEAN	Emergency access mode
<code>auth.jurisdiction()</code>	TEXT	User's legal jurisdiction
<code>auth.data_region()</code>	TEXT	Data residency region

## Context Management

```
-- Set context for request
SELECT auth.set_context(
  p_tenant_id := '...',
  p_user_id := '...',
  p_app_id := 'thinktank',
  p_permission_level := 'tenant_admin',
  p_jurisdiction := 'EU',
  p_data_region := 'eu-west-1',
  p_break_glass := false
);

-- Clear after request
SELECT auth.clear_context();
```

### 34.3 User Application Assignments

Users can be assigned to multiple applications with different permission levels.

#### Assignment Types

Type	Description
<b>standard</b>	Normal user access
<b>admin</b>	Application admin
<b>readonly</b>	Read-only access
<b>trial</b>	Limited trial access

#### API Endpoints

Endpoint	Method	Description
/user-registry/assignments	GET	List assignments (by userId or appId)
/user-registry/assignments	POST	Create assignment
/user-registry/assignments/revoke	POST	Revoke assignment

### 34.4 Consent Management

GDPR/CCPA/COPPA compliant consent tracking with full audit trail.

#### Lawful Bases (GDPR Article 6)

Basis	Description
<b>consent</b>	User explicitly consented
<b>contract</b>	Necessary for contract
<b>legal_obligation</b>	Required by law
<b>vital_interests</b>	Protect vital interests
<b>public_interest</b>	Public interest task
<b>legitimate_interests</b>	Legitimate business interest

## Consent Methods

Method	Description
<code>explicit_checkbox</code>	Unchecked checkbox
<code>click_accept</code>	Click to accept button
<code>implicit</code>	Implied consent
<code>parent_consent</code>	Parent/guardian consent
<code>verified_parent</code>	Verified parental consent (COPPA)
<code>double_opt_in</code>	Email confirmation

## API Endpoints

Endpoint	Method	Description
<code>/user-registry/consent</code>	GET	Get user consents
<code>/user-registry/consent</code>	POST	Record consent
<code>/user-registry/consent/withdraw</code>	POST	Withdraw consent
<code>/user-registry/consent/check</code>	GET	Check consent status

### 34.5 Break Glass Access

Emergency admin access for critical situations with full audit trail.

#### Requirements

- **Radiant Admin only** - Tenant admins cannot initiate
- **Incident ticket recommended** - Link to incident tracking
- **Approval tracking** - Record who approved access
- **P0 alerting** - SNS notification on initiation
- **Immutable audit log** - Hash-chained entries

## API Endpoints

Endpoint	Method	Description
<code>/user-registry/break-glass</code>	GET	List access logs
<code>/user-registry/break-glass/active</code>	GET	Active sessions
<code>/user-registry/break-glass/initiate</code>	POST	Start emergency access
<code>/user-registry/break-glass/end</code>	POST	End access session

### Example: Initiate Break Glass

POST `/api/admin/user-registry/break-glass/initiate`

```
{
  "tenantId": "tenant-uuid",
  "accessReason": "Customer escalation - data corruption investigation",
  "incidentTicket": "INC-2024-001234",
  "approvedBy": "cto@company.com"
}
```

Response:



```
{
  "success": true,
  "accessId": "bg-uuid",
  "message": "Break Glass access initiated",
  "instructions": "Access will be logged. End session when complete."
}
```

### 34.6 Legal Hold

Prevent data deletion for litigation, regulatory investigation, or audit.

#### API Endpoints

Endpoint	Method	Description
/user-registry/legal-hold	GET	List holds
/user-registry/legal-hold/apply	POST	Apply hold
/user-registry/legal-hold/release	POST	Release hold

#### Example: Apply Legal Hold

POST /api/admin/user-registry/legal-hold/apply

```
{
  "userId": "user-uuid",
  "reason": "Pending litigation - case #2024-CV-1234",
  "caseId": "2024-CV-1234"
}
```

### 34.7 DSAR Processing

Handle Data Subject Access Requests as required by GDPR, CCPA, etc.

#### Request Types

Type	Description	SLA
access	Export user data	30 days
delete	Delete user data	30 days
portability	Export in machine-readable format	30 days
rectification	Correct inaccurate data	30 days
restriction	Restrict processing	30 days
objection	Object to processing	30 days

#### API Endpoints

Endpoint	Method	Description
/user-registry/dsar	GET	List DSAR requests
/user-registry/dsar/process	POST	Process request
/user-registry/dsar/:id	PATCH	Update status

### 34.8 Cross-Border Transfer

Validate data transfers between regions based on user jurisdiction.

## Transfer Mechanisms

Mechanism	Description
same_region	Data stays in same region
pre_approved_region	Region in tenant's allowed list
adequacy_decision	EU adequacy decision exists
explicit_consent	User consented to transfer
sccs	Standard Contractual Clauses
bcr	Binding Corporate Rules

## API Endpoint

GET /api/admin/user-registry/cross-border/check?userId=xxx&targetRegion=us-west-2

### 34.9 Credential Rotation

Zero-downtime secret rotation for applications.

#### How It Works

1. **Generate new secret** - Create new credential
2. **Set rotation window** - Both secrets valid (default 24h)
3. **Update clients** - Migrate clients to new secret
4. **Window expires** - Old secret automatically invalid

#### API Endpoints

Endpoint	Method	Description
/user-registry/credentials/verify	POST	Verify credentials
/user-registry/credentials/rotate	POST	Start rotation
/user-registry/credentials/set	POST	Set initial secret
/user-registry/credentials/cleanup	POST	Clear expired windows

### 34.10 Infrastructure

#### DynamoDB Tables

Table	Purpose
radiant-app-client-mapping-{env}	M2M token enrichment
radiant-user-assignments-{env}	Assignment cache with TTL
radiant-tenant-config-{env}	Tenant configuration cache

#### S3 Audit Bucket

- **Object Lock:** COMPLIANCE mode, 7-year retention
- **Encryption:** CMK with key rotation
- **Lifecycle:** Glacier transition at 90 days
- **Partitioning:** year/month/day prefixes

**Security Alerts** SNS topic radiant-security-alerts-{env} for:

- Break Glass initiation
- Legal hold changes
- DSAR completions
- Suspicious access patterns

### 34.11 Database Migration

Migration file: packages/infrastructure/migrations/125\_multi\_app\_user\_registry.sql

Tables created/modified:

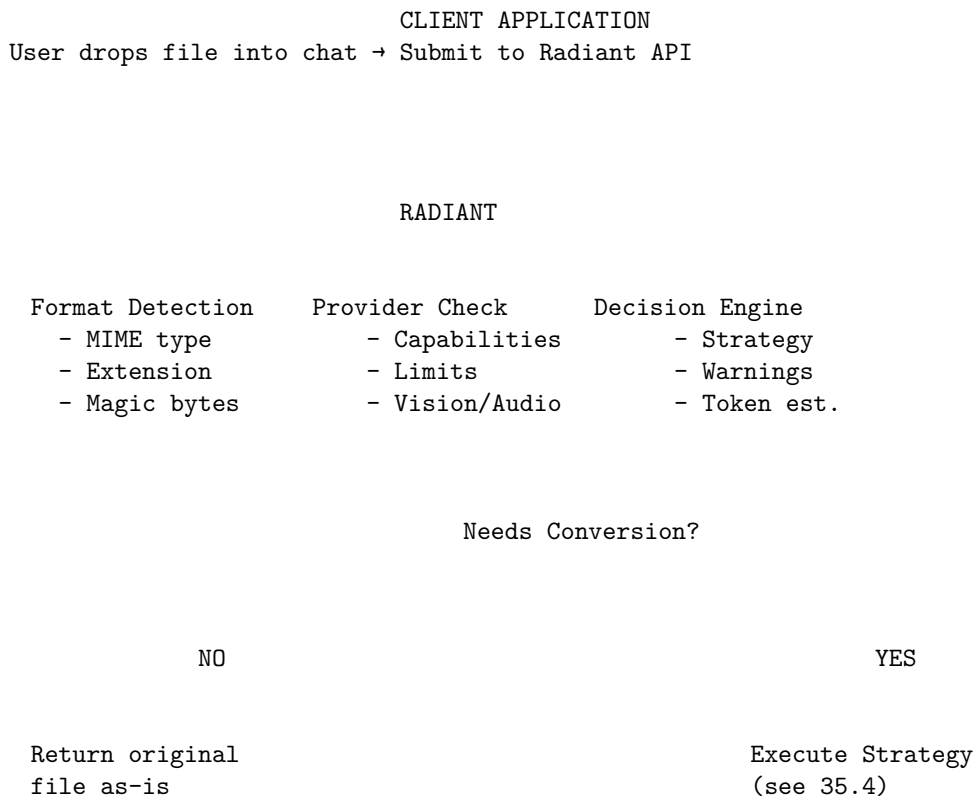
- Extended: tenants, users, registered\_apps
  - Created: user\_application\_assignments, consent\_records, data\_retention\_obligations, break\_glass\_access\_log, dsar\_requests
- 

## 35. Intelligent File Conversion Infrastructure

**Location:** Platform Infrastructure (operates automatically, no direct admin UI)

The **Intelligent File Conversion Service** is a Radiant-side system that automatically decides when and how to convert files for AI providers. Client applications (Think Tank, etc.) submit files; Radiant determines the optimal conversion strategy based on the target AI provider's capabilities.

### 35.1 Architecture Overview



Return Result  
+ Update Learning

## 35.2 Core Principle

”Let Radiant decide, not Think Tank”

1. Think Tank submits files **without worrying about provider compatibility**
2. Radiant detects file format and **checks target provider capabilities**
3. Conversion **only happens if the AI provider doesn't understand the format**
4. Uses **AI + libraries** for intelligent conversion
5. **Learns from outcomes** to improve future decisions

## 35.3 Provider Capabilities Registry

The service maintains a registry of what each AI provider can handle natively:

Provider	Vision	Audio	Video	Max File Size	Native Document Formats
<b>OpenAI</b>	GPT-4V	Whisper		20MB	txt, md, json, csv
<b>Anthropic</b>	Claude 3			32MB	pdf, txt, md, json, csv
<b>Google</b>	Gemini			100MB	pdf, txt, md, json, csv
<b>xAI</b>	Grok			20MB	txt, md, json
<b>DeepSeek</b>				10MB	txt, md, json, csv
<b>Self-hosted</b>	LLaVA	Whisper		50MB	txt, md, json, csv

## 35.4 Conversion Strategies (Detailed)

**none - No Conversion** Provider natively supports the format. File is passed through as-is.

**extract\_text - Text Extraction** Extracts plain text from documents.

Format	Library	Output
PDF	<b>pdf-parse</b>	Text + page metadata
DOCX/DOC	<b>mammoth</b>	Structured text
PPTX/PPT	native	Text from slides
HTML/XML	native	Tags stripped

**Example output:**

[Document Title]

Page 1:

Content from first page...

Page 2:

Content from second page...

[Metadata]

Pages: 10  
Author: John Doe  
Created: 2024-01-15

**ocr - Optical Character Recognition** Uses **AWS Textract** to extract text from images.

**Features:**

- Printed and handwritten text detection
- Table detection and extraction
- Form field detection
- Confidence scores per block

**Example output:**

[OCR Result]  
Confidence: 94.5%

INVOICE #12345  
Date: January 15, 2024

Item	Qty	Price
Widget A	10	\$50.00
Widget B	5	\$25.00

Total: \$625.00

**transcribe - Audio Transcription** Uses **OpenAI Whisper** (or self-hosted) for speech-to-text.

**Features:**

- Automatic language detection
- Timestamp segments
- SRT/VTT subtitle generation

**Example output:**

[Transcription]  
Duration: 5:32  
Language: English

[00:00] Hello and welcome to today's meeting.  
[00:05] We'll be discussing the Q4 roadmap.  
[00:12] First, let's review the current status...

**describe\_image - AI Image Description** Uses vision-capable models to describe image contents.

**Models used:**

- GPT-4 Vision (OpenAI)
- Claude 3 Vision (Anthropic)
- LLaVA (self-hosted)

**Example output:**

[Image Description]  
Model: gpt-4-vision  
Dimensions: 1920x1080

This image shows a modern office space with an open floor plan. In the foreground, there are several desks arranged in clusters, each with monitors and office supplies.

[Text detected in image]:  
"RADIANT - Innovation Center"

**describe\_video - Video Frame Analysis** Extracts key frames and describes each using vision models.

**Configuration:**

- Frame interval: 10 seconds (default)
- Maximum frames: 10 (default)

**Example output:**

**\*\*Video Overview\*\*** (2m 30s, 1920x1080)

**\*\*Frame Analysis:\*\***

**\*\*[0:00]\*\*** The video opens with a title screen showing the company logo.

**\*\*[0:10]\*\*** A presenter stands in front of a whiteboard with diagrams.

**\*\*[0:20]\*\*** Close-up of the whiteboard showing a flowchart.

**\*\*Summary:\*\***

The video begins with company logo and ends with key point summary.

**parse\_data - Structured Data Parsing**

Format	Library	Output
CSV	native	JSON array of objects
XLSX/XLS	<code>xlsx</code>	JSON with sheet data
JSON	native	Validated and prettified

**Example output (CSV):**

```
{
  "data": [
    {"name": "Alice", "email": "alice@example.com", "role": "Admin"},
    {"name": "Bob", "email": "bob@example.com", "role": "User"}
  ],
  "metadata": {
    "rowCount": 2,
    "columnCount": 3,
    "headers": ["name", "email", "role"]
  }
}
```

**decompress - Archive Extraction** Extracts and processes archive contents.

Format	Library
ZIP	<code>adm-zip</code>
TAR	<code>tar</code>
GZIP	<code>zlib</code>

Format	Library
--------	---------

#### Features:

- Recursive extraction
- Text file content inclusion
- Binary file detection
- Size limits enforcement

**render\_code - Code Formatting** Formats code files with syntax highlighting as markdown.

### 35.5 Supported File Formats

#### Documents

Format	Extension	Conversion Strategy
PDF	.pdf	<code>extract_text</code> via pdf-parse
Word	.docx, .doc	<code>extract_text</code> via mammoth
PowerPoint	.pptx, .ppt	<code>extract_text</code>
Excel	.xlsx, .xls	<code>parse_data</code> via xlsx

#### Text Files

Format	Extension	Notes
Plain Text	.txt	Direct passthrough
Markdown	.md	Direct passthrough
JSON	.json	Direct or <code>parse_data</code>
CSV	.csv	<code>parse_data</code>
XML	.xml	Direct or <code>extract_text</code>
HTML	.html	<code>extract_text</code>

#### Images

Format	Extension	Conversion Strategy
PNG	.png	Native or <code>describe_image</code>
JPEG	.jpg, .jpeg	Native or <code>describe_image</code>
GIF	.gif	Native or <code>describe_image</code>
WebP	.webp	Native or <code>describe_image</code>
SVG	.svg	Convert to PNG or <code>describe_image</code>
BMP	.bmp	Convert to PNG or <code>describe_image</code>
TIFF	.tiff	Convert to PNG or <code>describe_image</code>

#### Audio

Format	Extension	Conversion Strategy
MP3	.mp3	<code>transcribe</code> via Whisper
WAV	.wav	<code>transcribe</code> via Whisper
OGG	.ogg	<code>transcribe</code> via Whisper

Format	Extension	Conversion Strategy
FLAC	.flac	transcribe via Whisper
M4A	.m4a	transcribe via Whisper

## Video

Format	Extension	Conversion Strategy
MP4	.mp4	describe_video
WebM	.webm	describe_video
MOV	.mov	describe_video
AVI	.avi	describe_video

## Code Files

Format	Extension	Notes
Python	.py	Syntax-highlighted markdown
JavaScript	.js, .jsx	Syntax-highlighted markdown
TypeScript	.ts, .tsx	Syntax-highlighted markdown
Java	.java	Syntax-highlighted markdown
C/C++	.c, .cpp, .h	Syntax-highlighted markdown
Go	.go	Syntax-highlighted markdown
Rust	.rs	Syntax-highlighted markdown
Ruby	.rb	Syntax-highlighted markdown

## Archives

Format	Extension	Conversion Strategy
ZIP	.zip	decompress via adm-zip
TAR	.tar	decompress via tar
GZIP	.gz, .tar.gz	decompress via zlib

## 35.6 Multi-Model File Preparation

When multiple AI models work on the same prompt (multi-model orchestration), the system makes **per-model conversion decisions**:

**Key Principle:** "If a model accepts the file type, assume it understands it unless proven otherwise."

File: document.pdf → 3 Models

Claude 3.5	GPT-4	DeepSeek
PDF:	PDF:	PDF:
PASS	CONVERT	CONVERT
ORIGINAL	(extract)	(cached)



### Per-Model Actions:

Action	When	Result
pass_original	Model natively supports format	Original file passed
convert	Model doesn't support format	Converted content passed
skip	File too large or conversion failed	Model excluded

### Features:

- **Cached conversions:** Convert once, reuse for all models that need it
- **Per-model capability checking:** Vision, audio, video, document formats
- **Learned understanding:** Uses reinforcement learning data (see 35.10)

## 35.7 Domain-Specific File Formats

The service includes a registry of **50+ domain-specific formats** that are widely used in specialized fields but not commonly supported by mainstream AI providers.

### Mechanical Engineering / CAD

Format	Extensions	Description	Recommended Library
<b>STEP</b>	.step, .stp, .p21	ISO 10303 CAD exchange	OpenCASCADE, FreeCAD
<b>STL</b>	.stl	3D printing mesh	numpy-stl, trimesh
<b>OBJ</b>	.obj	Wavefront 3D model	trimesh, three.js
<b>Fusion 360</b>	.f3d, .f3z	Autodesk parametric CAD	Fusion 360 API
<b>IGES</b>	.iges, .igs	Legacy CAD exchange	OpenCASCADE
<b>DXF</b>	.dxf	AutoCAD 2D drawings	ezdxf
<b>GLTF/GLB</b>	.gltf, .glb	Web 3D format	three.js, trimesh

### CAD Converter Extracts:

Format	What's Extracted
<b>STL</b>	Triangle count, bounding box, 3D printing assessment, volume estimate
<b>OBJ</b>	Vertices, faces, materials, groups
<b>STEP</b>	Entities, part names, assembly structure, schema version
<b>DXF</b>	Layers, entity types (lines, arcs, circles), block count
<b>GLTF/GLB</b>	Meshes, materials, animations, scene graph

### Electrical Engineering

Format	Extensions	Description	Library
<b>KiCad</b>	.kicad_pcb, .kicad_sch	PCB/schematic	kicad-cli, kiutils
<b>EAGLE</b>	.brd, .sch	Autodesk PCB	eagle-to-kicad
<b>SPICE</b>	.spice, .sp, .cir	Circuit simulation	PySpice, ngspice

### Medical/Healthcare

Format	Extensions	Description	Library
<b>DICOM</b>	.dcm, .dicom	Medical imaging	pydicom, dcmtk
<b>HL7 FHIR</b>	.json, .xml	Health records	fhir.resources

## Scientific/Research

Format	Extensions	Description	Library
<b>NetCDF</b>	.nc, .nc4	Climate/geoscience	netCDF4, xarray
<b>HDF5</b>	.h5, .hdf5	Scientific data	h5py
<b>FITS</b>	.fits	Astronomy data	astropy

## Geospatial

Format	Extensions	Description	Library
<b>Shapefile</b>	.shp, .dbf	Vector GIS	geopandas, shapefile
<b>GeoTIFF</b>	.tif, .geotiff	Georeferenced raster	rasterio

## Bioinformatics

Format	Extensions	Description	Library
<b>FASTA</b>	.fasta, .fa	DNA/protein sequences	Biopython
<b>PDB</b>	.pdb	Protein structure	Biopython, py3Dmol

**Domain Detection** When a domain-specific file is uploaded:

1. Format detected by extension and MIME type
2. Domain identified (e.g., Mechanical Engineering)
3. AGI Brain selects appropriate library
4. Extracts domain-relevant information
5. Generates AI-readable description with specialized prompts

### Example AI Prompts per Format:

```
# STL file prompt
"This is an STL 3D model file. Describe the shape, identify what object
it might be, assess printability, and note any potential issues for 3D printing."
```

```
# DICOM file prompt
"This is a DICOM medical image. Describe the imaging modality, anatomical
region, and any visible findings. Note: Do not provide medical diagnoses."
```

```
# STEP file prompt
"This is a STEP CAD file. Describe the mechanical part or assembly,
including approximate geometry, features, and likely manufacturing process."
```

## 35.8 Database Schema

Migration: 127\_file\_conversion\_service.sql

Table	Purpose
<code>file_conversions</code>	Tracks all conversion decisions and results
<code>provider_file_capabilities</code>	Provider format support registry (configurable)

#### `file_conversions` columns:

- `id` - UUID primary key
- `tenant_id` - Tenant reference
- `user_id` - User who uploaded
- `conversation_id` - Associated conversation
- `filename` - Original filename
- `mime_type` - MIME type
- `file_format` - Detected format
- `file_size` - Size in bytes
- `file_checksum` - SHA256 hash
- `target_provider_id` - Target AI provider
- `target_model_id` - Target model
- `conversion_strategy` - Strategy used
- `conversion_decision` - Full decision JSON
- `converted_content` - Converted content (if applicable)
- `token_estimate` - Estimated tokens
- `processing_time_ms` - Processing duration
- `success` - Success boolean
- `error` - Error message if failed
- `created_at` - Timestamp

**View:** `v_file_conversion_stats` - Aggregated statistics per tenant

#### Functions:

- `check_format_supported(provider, format)` - Check if format is natively supported
- `get_conversion_stats(tenant_id, days)` - Get tenant statistics for N days
- `cleanup_old_conversions(retention_days)` - Clean up old conversion records

**Migration:** `128_file_conversion_learning.sql`

Table	Purpose
<code>model_format_understanding</code>	Per-tenant learned model/format scores
<code>conversion_outcome_feedback</code>	Recorded feedback for learning
<code>format_understanding_events</code>	Audit trail of score changes
<code>global_format_learning</code>	Cross-tenant aggregate insights

## 35.9 API Endpoints

**Base Path:** `/api/thinktank/files`

### Process File

POST `/api/thinktank/files/process`

#### Request:

```
{
  "filename": "document.pdf",
  "mimeType": "application/pdf",
```

```

    "content": "<base64-encoded-content>",
    "targetProvider": "anthropic",
    "targetModel": "claude-3-5-sonnet",
    "conversationId": "conv-uuid-optional"
}

```

#### Response:

```

{
  "success": true,
  "data": {
    "conversionId": "conv_abc123",
    "originalFile": {
      "filename": "document.pdf",
      "format": "pdf",
      "size": 1048576,
      "checksum": "sha256:abc123..."
    },
    "convertedContent": {
      "type": "text",
      "content": "Extracted document text...",
      "tokenEstimate": 2500,
      "metadata": {
        "originalFormat": "pdf",
        "conversionStrategy": "extract_text",
        "pageCount": 10,
        "title": "Annual Report 2024"
      }
    },
    "processingTimeMs": 1250
  }
}

```

#### Check Compatibility

POST /api/thinktank/files/check-compatibility

#### Request:

```

{
  "filename": "image.png",
  "mimeType": "image/png",
  "fileSize": 524288,
  "targetProvider": "deepseek"
}

```

#### Response:

```

{
  "success": true,
  "data": {
    "fileInfo": {
      "filename": "image.png",
      "format": "png",
      "size": 524288
    },
    "provider": {

```

```

    "id": "deepseek",
    "supportsFormat": false,
    "supportsVision": false,
    "maxFileSize": 10485760
  },
  "decision": {
    "needsConversion": true,
    "strategy": "describe_image",
    "reason": "Provider deepseek lacks vision - will use AI to describe image",
    "targetFormat": "txt"
  }
}
}

```

### Get Capabilities

GET /api/thinktank/files/capabilities

Returns all provider capabilities for UI display.

### Get History

GET /api/thinktank/files/history?limit=50&offset=0

Returns conversion history for the current user.

### Get Statistics

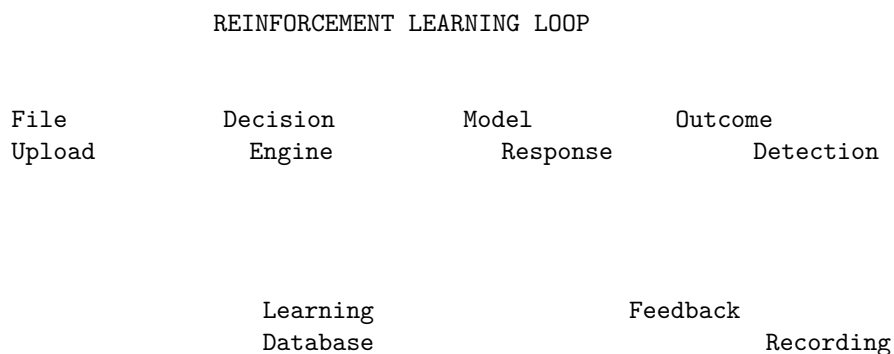
GET /api/thinktank/files/stats?days=30

Returns conversion statistics for the current tenant.

## 35.10 Reinforcement Learning Integration

The file conversion system integrates with the AGI Brain/consciousness for **persistent learning from conversion outcomes**.

### How Learning Works



### Learning Signals

Signal	Source	What It Learns
<b>User Rating</b>	Explicit 1-5 stars	Direct quality signal
<b>Model Response</b>	Auto-inferred from text	Did model understand?
<b>Error Detection</b>	Model errors/hallucinations	Format incompatibility
<b>Conversion Outcome</b>	Success/failure	Model capability

**Understanding Score** Each model/format combination has a learned understanding score (0.0 to 1.0):

Score	Level	Recommended Action
0.8 - 1.0	Excellent	Pass original file
0.6 - 0.8	Good	Pass original file
0.4 - 0.6	Moderate	May need conversion
0.0 - 0.4	Poor	Always convert

**Auto-Inference from Response** The system automatically detects outcomes from model responses:

**Failure signals detected:**

- "I can't read", "unable to process", "cannot access the file"
- "appears to be empty", "binary data", "base64"
- Model asking for clarification about file content
- Hallucinated content detection

**Consciousness Integration** Significant learning events create **Learning Candidates** for the consciousness system:

Event	Learning Candidate Type	Quality
Model failed on format it claimed to support	<code>format_misunderstanding</code>	0.85
Unnecessary conversion (model would have understood)	<code>unnecessary_conversion</code>	0.70
Model hallucinated file content	<code>hallucination_detection</code>	0.90
User gave negative rating	<code>user_correction</code>	0.85

These feed into the **LoRA evolution system** for persistent consciousness improvement.

**Admin Override** Force conversion for problematic model/format combinations:

```
-- Check current understanding scores
SELECT model_id, file_format, understanding_score, confidence,
       success_count, failure_count
FROM model_format_understanding
WHERE tenant_id = 'your-tenant-id'
ORDER BY understanding_score ASC;
```

**API Override:**

```
POST /api/admin/file-conversion/force-convert
{
  "modelId": "claude-3-haiku",
  "fileFormat": "pdf",
  "reason": "Struggles with multi-column PDFs"
}
```

```
DELETE /api/admin/file-conversion/force-convert
{
  "modelId": "claude-3-haiku",
  "fileFormat": "pdf"
}
```

### 35.11 Environment Variables

Variable	Description	Default	Required
FILE_CONVERSION_BUCKET	S3 bucket for file storage	radiant-files	Yes
OPENAI_API_KEY	OpenAI API key for Whisper/Vision	-	Yes
ANTHROPIC_API_KEY	Anthropic API key for Claude Vision	-	No
WHISPER_ENDPOINT_URL	Self-hosted Whisper endpoint	-	No
VISION_ENDPOINT_URL	Self-hosted LLaVA endpoint	-	No
AWS_TEXTTRACT_REGION	AWS region for Textract	us-east-1	No
FILE_CONVERSION_MAX_SIZE_MB	Maximum file size	100	No
FILE_CONVERSION_TIMEOUT_SEC	Processing timeout	60	No

### 35.12 Implementation Files

File	Purpose
lambda/shared/services/file-conversion.service.ts	Main conversion service with decision engine
lambda/shared/services/multi-model-file-prep.service.ts	Multi-model preparation
lambda/shared/services/file-conversion-learning.service.ts	Reinforcement learning
lambda/shared/services/converters/pdf-converter.ts	PDF text extraction
lambda/shared/services/converters/docx-converter.ts	DOCX extraction
lambda/shared/services/converters/excel-converter.ts	Excel/CSV parsing
lambda/shared/services/converters/audio-converter.ts	Audio transcription
lambda/shared/services/converters/image-converter.ts	Image description + OCR
lambda/shared/services/converters/video-converter.ts	Video frame extraction
lambda/shared/services/converters/archive-converter.ts	Archive decompression
lambda/shared/services/converters/cad-converter.ts	CAD/3D file parsing
lambda/shared/services/converters/domain-formats.ts	Domain format registry
lambda/shared/services/converters/domain-converter-selector.ts	AGI Brain integration
lambda/thinktank/file-conversion.ts	API handlers
migrations/127_file_conversion_service.sql	Main database schema
migrations/128_file_conversion_learning.sql	Learning schema

### 35.13 Monitoring

#### Metrics tracked per tenant:

- Total files processed
- Conversion success/failure rate
- Average processing time
- Most common formats
- Most common conversion strategies
- Storage usage
- Learning statistics (formats learned, understanding improvements)

#### Alerts:

- High failure rate (>5%)

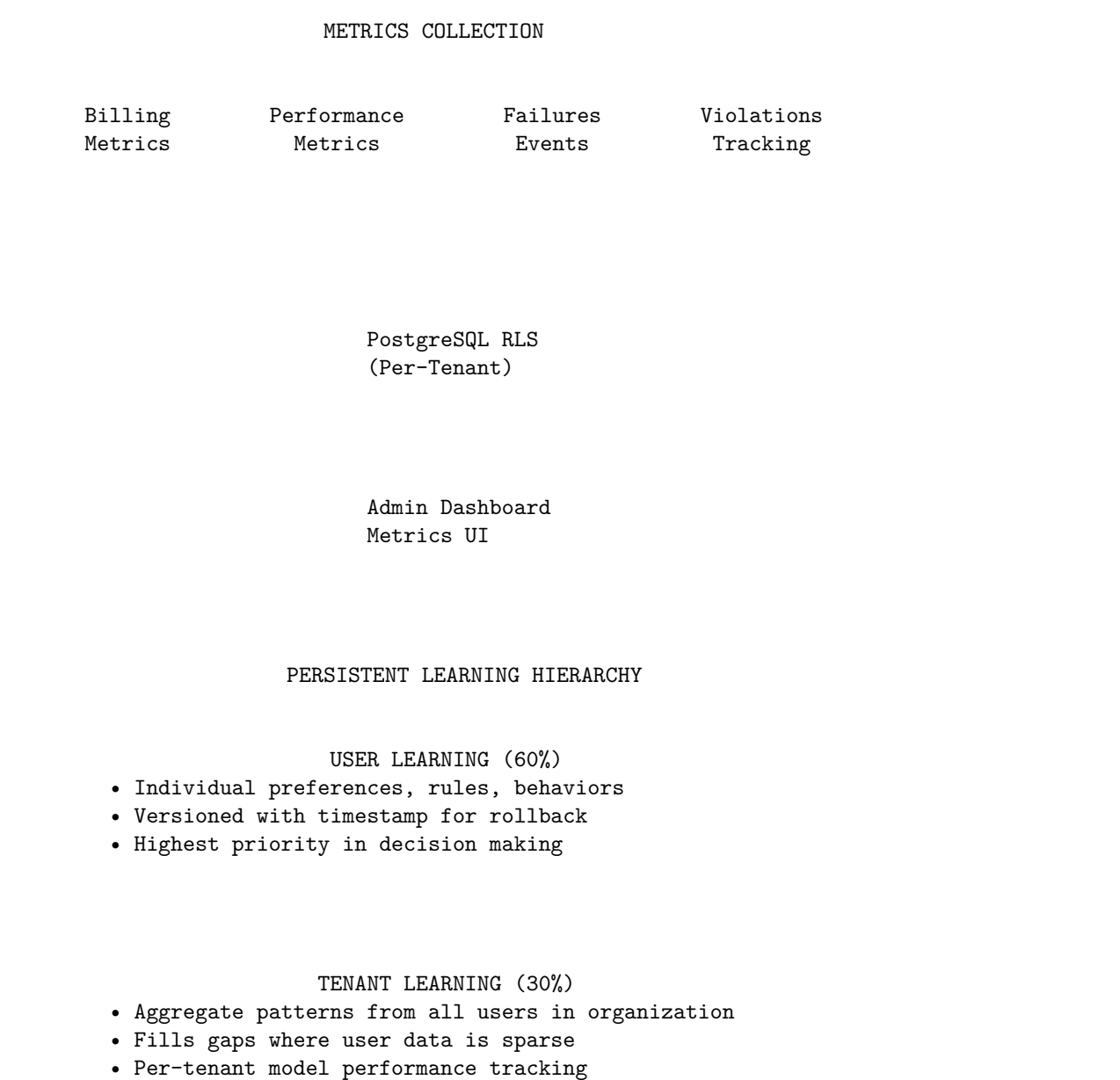
- Processing time > 30s
  - Storage quota approaching limit
  - Learning anomalies (sudden score drops)
- 

## 36. Metrics & Persistent Learning Infrastructure

**Location:** Admin Dashboard → Metrics

The Metrics & Persistent Learning Infrastructure provides comprehensive tracking of billing, performance, failures, violations, and logs, plus a persistent learning system that survives system reboots with a User → Tenant → Global influence hierarchy.

### 36.1 Architecture Overview





#### GLOBAL LEARNING (10%)

- Anonymized aggregate from all tenants
- Baseline intelligence, privacy-protected
- Minimum 5 tenant threshold for aggregation

#### SNAPSHOT & RECOVERY

- Daily snapshots for fast recovery
- System reboots do NOT require relearning
- Checksums for integrity verification

### 36.2 Metrics Categories

**Billing Metrics** Tracks cost and usage per tenant/user/model:

- **Token usage:** Input/output tokens per request
- **Cost breakdown:** Cost in cents (avoids floating point issues)
- **Storage usage:** Bytes used and storage cost
- **Compute usage:** Self-hosted compute seconds and cost
- **API call counts:** Successful and failed calls

**Performance Metrics** Tracks latency and throughput:

- **Total latency:** End-to-end request time
- **Time to first token:** Streaming response start
- **Inference time:** Model processing time
- **Queue wait time:** Time in request queue
- **Throughput:** Tokens per second

**Failure Events** Tracks errors with classification:

Failure Type	Description
api_error	API endpoint errors
model_error	Model inference failures
timeout	Request timeouts
rate_limit	Rate limiting triggered
auth_error	Authentication failures
validation_error	Request validation failures
provider_error	External provider failures
quota_exceeded	Usage quota exceeded
content_filter	Content filtered
context_length	Context too long

Severity levels: low, medium, high, critical

**Prompt Violations** Tracks content policy violations:

Violation Type	Description
content_policy	General policy violation
jailbreak_attempt	Attempted jailbreak
injection_attempt	Prompt injection
pii_exposure	PII in prompt
harassment	Harassment content
hate_speech	Hate speech detected
violence	Violence content
sexual_content	Sexual content
illegal_activity	Illegal activity

**Actions taken:** blocked, warned, filtered, logged, escalated

**System Logs** Centralized logging with levels:

- trace, debug, info, warn, error, fatal

### 36.3 User Learning (Think Tank Rules Included)

**User Rules (Versioned)** User-defined rules for AI behavior with automatic versioning:

```
interface UserRule {
  id: string;
  tenantId: string;
  userId: string;
  ruleName: string;
  ruleCategory: 'behavior' | 'format' | 'tone' | 'content' | 'restriction' |
    'preference' | 'domain' | 'persona' | 'workflow' | 'other';
  currentVersion: number;
  ruleContent: string;
  rulePriority: number; // 0-100
  isActive: boolean;
  appliesToModels?: string[];
  appliesToTasks?: string[];
  effectivenessScore: number; // Learned from outcomes
}
```

**Automatic versioning:** Every update creates a new version for rollback capability.

**User Learned Preferences** AGI Brain learns user preferences automatically:

Category	Examples
communication_style	Formal, casual, technical
response_format	Bullet points, paragraphs, code
detail_level	Brief, detailed, comprehensive
expertise_level	Beginner, intermediate, expert
model_preference	Preferred models for tasks
language	Response language preference

**Learning sources:**

- `explicit_setting` - User explicitly set
- `implicit_behavior` - Inferred from behavior
- `feedback` - Derived from ratings
- `conversation` - Learned from chat
- `pattern_detection` - Automatic pattern matching

### 36.4 Tenant Aggregate Learning

Aggregates learning across all users in a tenant:

#### Learning Dimensions

- **Model performance:** Quality, speed, cost-efficiency, reliability per model/task
- **Task patterns:** Common task types and successful approaches
- **Error recovery:** How to recover from specific errors
- **Format preferences:** Organization-wide format preferences
- **Domain expertise:** Learned domain-specific knowledge

**Model Performance Tracking** Per-tenant model scoring:

```
SELECT model_id, task_type, quality_score, reliability_score,
       total_uses, successful_uses, confidence
FROM tenant_model_performance
WHERE tenant_id = 'your-tenant-id'
ORDER BY quality_score DESC;
```

### 36.5 Global Aggregate Learning

Anonymized cross-tenant learning:

- **Minimum threshold:** 5 tenants before data is included
- **Anonymization:** User/tenant data is hashed, never stored
- **Privacy controls:** Per-tenant opt-out available
- **Pattern library:** Successful patterns shared anonymously

### 36.6 Learning Influence Configuration

Per-tenant configuration for influence weights:

```
interface LearningInfluenceConfig {
    tenantId: string;
    userWeight: number;    // Default: 0.60
    tenantWeight: number;  // Default: 0.30
    globalWeight: number;  // Default: 0.10
    enableUserLearning: boolean;
    enableTenantAggregation: boolean;
    enableGlobalLearning: boolean;
    contributeToGlobal: boolean;
}
```

Weights must sum to 1.0

### 36.7 Persistence & Recovery

**Snapshots** Daily snapshots for fast recovery:

- **User snapshots:** All preferences, rules, memories
- **Tenant snapshots:** Aggregate learning, model performance

- **Global snapshots:** Cross-tenant aggregates, pattern library

**Recovery Process** On system reboot or failure:

1. Load latest snapshot
2. Verify checksum integrity
3. Restore learning state
4. Log recovery event

**NO RELEARNING REQUIRED** - All learning is persisted in PostgreSQL.

### 36.8 Database Schema

Migration: 129\_metrics\_persistent\_learning.sql

#### Metrics Tables

Table	Purpose
billing_metrics	Cost and usage tracking
performance_metrics	Latency and throughput
failure_events	Error tracking
prompt_violations	Policy violations
system_logs	Centralized logs
mv_tenant_daily_metrics	Materialized view for daily aggregates

#### Learning Tables

Table	Purpose
user_rules	User-defined rules (versioned)
user_rules_versions	Rule version history
user_learned_preferences	Learned user preferences
user_preference_versions	Preference version history
user_memory_versions	Memory version history
tenant_aggregate_learning	Tenant-level learning state
tenant_learning_events	Tenant learning audit trail
tenant_model_performance	Per-tenant model scores
global_aggregate_learning	Cross-tenant learning
global_model_performance	Global model scores
global_pattern_library	Shared patterns
learning_influence_config	Per-tenant influence weights
learning_decision_log	Decision audit trail
learning_snapshots	Point-in-time snapshots
learning_recovery_log	Recovery audit trail

### 36.9 API Endpoints

Base Path: /api/admin/metrics

#### Dashboard & Summary

Endpoint	Method	Description
/dashboard	GET	Full dashboard data
/summary	GET	Metrics summary

## Billing

Endpoint	Method	Description
/billing	GET	Billing metrics
/billing	POST	Record billing metric

## Performance

Endpoint	Method	Description
/performance	GET	Performance metrics
/performance/latency	GET	Latency percentiles
/performance	POST	Record performance metric

## Failures

Endpoint	Method	Description
/failures	GET	Failure events
/failures	POST	Record failure
/failures/:id/resolve	POST	Resolve failure

## Violations

Endpoint	Method	Description
/violations	GET	Prompt violations
/violations	POST	Record violation
/violations/:id/review	POST	Review violation

## Learning

Endpoint	Method	Description
/learning/influence	GET	Get learning influence
/learning/config	GET/PUT	Influence configuration
/learning/tenant	GET	Tenant learning state
/learning/global	GET	Global learning state
/learning/model-performance	GET	Model performance scores
/learning/event	POST	Record learning event
/learning/user-preferences	GET/POST	User preferences

## Snapshots & Recovery

Endpoint	Method	Description
/learning/snapshots	GET	Get latest snapshot
/learning/snapshots	POST	Create snapshot
/learning/snapshots/:id/recover	POST	Recover from snapshot
/learning/recovery-logs	GET	Recovery history

### 36.10 Admin Dashboard

Location: apps/admin-dashboard/app/(dashboard)/metrics/page.tsx

#### Dashboard Tabs

- **Overview:** Summary cards, daily usage chart, top models
- **Billing:** Cost breakdown by date and model
- **Performance:** Latency percentiles, model performance table
- **Failures:** Recent failures with severity and resolution status
- **Violations:** Content violations with review status
- **Learning:** Learning hierarchy visualization, snapshot status

### 36.11 Implementation Files

File	Purpose
packages/shared/src/types/metrics-learning.types.ts	TypeScript types
lambda/shared/services/metrics-collection.service.ts	Metrics collection
lambda/shared/services/learning-influence.service.ts	Learning hierarchy
lambda/shared/middleware/metrics-middleware.ts	Auto-metrics for AI endpoints
lambda/admin/metrics.ts	API handlers
lambda/scheduled/learning-snapshots.ts	Daily snapshot Lambda
lambda/scheduled/learning-aggregation.ts	Weekly aggregation Lambda
migrations/129_metrics_persistent_learning.sql	Database schema
apps/admin-dashboard/app/(dashboard)/metrics/page.tsx	Admin UI
lib/stacks/api-stack.ts	CDK routes (lines 463-625)
lib/stacks/scheduled-tasks-stack.ts	Scheduled Lambdas

### 36.15 Integration Points

**Cognitive Router Integration** The Cognitive Router automatically uses learning influence when `useLearningInfluence: true`:

```
import { brainRouter, initializeLearningService } from './services/cognitive-router.service';

// Initialize once at startup
initializeLearningService(pool);

// Route with learning influence
const result = await brainRouter.route({
  tenantId,
  userId,
  taskType: 'code',
  prompt: userPrompt,
  useLearningInfluence: true, // Enable User → Tenant → Global
  useDomainProficiencies: true,
  useAffectMapping: true,
```

```
});

// result.learningDecisionId - Use for feedback loop
// result.learningInfluenceUsed - true if learning was applied

// Record outcome after user feedback
await brainRouter.recordRoutingOutcome(tenantId, result.learningDecisionId, positive);
```

**Metrics Middleware** Wrap AI endpoints to automatically record metrics:

```
import { withMetrics } from './middleware/metrics-middleware';

// Wrap handler
export const handler = withMetrics(async (event, context) => {
  // Your handler logic
  return { statusCode: 200, body: JSON.stringify(result) };
});
```

**Manual Metrics Recording** Record metrics programmatically:

```
import {
  recordBillingMetric,
  recordFailure,
  recordViolation,
  logSystem
} from './middleware/metrics-middleware';

// Record billing
await recordBillingMetric(tenantId, userId, modelId, inputTokens, outputTokens, costCents);

// Record failure
await recordFailure(tenantId, userId, 'model_error', 'high', 'Model timeout', modelId);

// Record violation
await recordViolation(tenantId, userId, 'jailbreak_attempt', 'high', promptSnippet, 'blocked');

// System log
await logSystem('info', 'my-service', 'Operation completed', { details: 'value' }, tenantId);
```

### 36.16 Scheduled Tasks

Task	Schedule	Purpose
<b>Learning Snapshots</b>	Daily 3 AM UTC	Create user/tenant/global learning backups
<b>Learning Aggregation</b>	Weekly Sun 4 AM UTC	Aggregate tenant data to global (min 5 tenants)

Both Lambdas are configured in `scheduled-tasks-stack.ts` and handle:

- Error recovery with logging
- Cleanup of old data (30-day snapshots, 90-day events)
- Materialized view refresh for dashboard performance

### 36.12 Environment Variables

Variable	Description	Default
METRICS_RETENTION_DAYS	Days to retain detailed metrics	90
SNAPSHOT_FREQUENCY_HOURS	Hours between snapshots	24
GLOBAL_AGGREGATION_MIN_TENANTS	Minimum tenants for global	5
LEARNING_WEIGHTS_USER	Default user weight	0.60
LEARNING_WEIGHTS_TENANT	Default tenant weight	0.30
LEARNING_WEIGHTS_GLOBAL	Default global weight	0.10

### 36.13 Monitoring & Alerts

#### Metrics to monitor:

- Billing metrics recording rate
- Failure rate by severity
- Violation rate by type
- Snapshot success rate
- Recovery success rate
- Learning event volume

#### Alerts:

- High failure rate (>5%)
- Critical severity failures
- Snapshot failures
- Recovery failures
- Unusual violation patterns

### 36.14 Security Considerations

- **RLS enforcement:** All tables use Row Level Security
- **Super admin bypass:** For cross-tenant analytics
- **Anonymization:** Global learning uses hashed identifiers
- **Audit trail:** All decisions logged for compliance
- **Data retention:** Configurable retention periods

---

## 37. Translation Middleware (18 Language Support)

The Translation Middleware provides automatic translation for multilingual support across all 18 supported languages. It intelligently routes prompts through translation when the target model doesn't natively support the user's language.

### 37.1 Overview

**Purpose:** Enable cost-effective self-hosted models to serve users in any of the 18 supported languages by transparently translating input to English, processing, and translating output back.

#### Architecture:

User Input (Any Language)

Language Detection



Model Language Matrix Check  
(Is translation required?)

No

Yes

Direct  
Process

Translate to  
English

Process with  
Model

Translate  
Back

User Output

### 37.2 Supported Languages

#	Code	Language	Native Name	RTL
1	en	English	English	No
2	es	Spanish	Español	No
3	fr	French	Français	No
4	de	German	Deutsch	No
5	pt	Portuguese	Português	No
6	it	Italian	Italiano	No
7	nl	Dutch	Nederlands	No
8	pl	Polish	Polski	No
9	ru	Russian		No
10	tr	Turkish	Türkçe	No
11	ja	Japanese		No
12	ko	Korean		No
13	zh-CN	Chinese (Simplified)		No
14	zh-TW	Chinese (Traditional)		No
15	ar	Arabic		<b>Yes</b>
16	hi	Hindi		No
17	th	Thai		No
18	vi	Vietnamese	Tiếng Việt	No

### 37.3 Model Language Support Levels

Each model has a language capability matrix defining support levels:

Level	Quality Score	Behavior
native	90-100	No translation needed
good	75-89	No translation (acceptable quality)
moderate	50-74	May translate depending on threshold
poor	25-49	Translation recommended
none	0-24	Translation required

Model Categories:

Model Type	Translate Threshold	Example
External (Claude, GPT-4)	<b>none</b>	Never translate
Large Self-Hosted (Qwen 72B)	<b>moderate</b>	Translate for poor/none
Medium Self-Hosted (Llama 70B)	<b>good</b>	Translate for moderate/poor/none
Small Self-Hosted (7B models)	<b>native</b>	Translate for all non-English

### 37.4 Translation Model

Default: qwen2.5-7b-instruct

Why Qwen 2.5 7B:

- Excellent multilingual capabilities (trained on Chinese + English + 27 languages)
- Cost-effective: \$0.08/1M input, \$0.24/1M output
- 3x cheaper than Claude Haiku
- Fast inference on g5.2xlarge

**Cost Comparison:** | Model | Input/1M | Output/1M | |-----|-----|-----| | **Qwen 2.5 7B** | \$0.08 | \$0.24 | | Claude Haiku | \$0.25 | \$1.25 | | GPT-3.5 Turbo | \$0.50 | \$1.50 |

### 37.5 Configuration

Admin UI: Settings → Translation

Configuration Options:

Setting	Default	Description
enabled	true	Enable/disable translation middleware
translation_model	qwen2.5-7b-instruct	Model used for translation
cache_enabled	true	Cache translations to reduce cost
cache_ttl_hours	168	Cache expiry (7 days)
max_cache_size	10000	Max cache entries per tenant
confidence_threshold	0.70	Min confidence to accept translation
max_input_length	50000	Max characters per translation
preserve_code_blocks	true	Don't translate code
preserve_urls	true	Don't translate URLs
preserve_mentions	true	Don't translate @mentions
fallback_to_english	true	Fallback if translation fails
cost_limit_per_day_cents	1000	Max daily translation cost (\$10)

## 37.6 API Reference

Base URL: /api/admin/translation

### Get Configuration

GET /api/admin/translation/config

Response:

```
{
  "config": {
    "enabled": true,
    "translation_model": "qwen2.5-7b-instruct",
    "cache_enabled": true,
    "cache_ttl_hours": 168,
    "cost_limit_per_day_cents": 1000
  },
  "isDefault": false
}
```

### Update Configuration

PUT /api/admin/translation/config

Content-Type: application/json

```
{
  "enabled": true,
  "cost_limit_per_day_cents": 2000
}
```

### Get Dashboard

GET /api/admin/translation/dashboard?days=30

Response:

```
{
  "config": { ... },
  "metrics": {
    "totalTranslations": 1523,
    "byLanguagePair": {
      "ja->en": 450,
      "en->ja": 448,
      "zh-CN->en": 312,
      "en->zh-CN": 313
    },
    "totalTokens": { "input": 2500000, "output": 2800000 },
    "estimatedCost": 0.87,
    "periodDays": 30
  },
  "cache": {
    "entriesCount": 856,
    "totalHits": 4521
  },
  "recentTranslations": [ ... ]
}
```

## Detect Language

POST /api/admin/translation/detect  
Content-Type: application/json

```
{  
  "text": "  
}
```

### Response:

```
{  
  "detectedLanguage": "ja",  
  "confidence": 0.95,  
  "alternativeLanguages": [  
    { "language": "zh-CN", "confidence": 0.03 }  
  ],  
  "isMultilingual": false,  
  "scriptType": "cjk"  
}
```

## Check Model Language Support

POST /api/admin/translation/check-model  
Content-Type: application/json

```
{  
  "modelId": "llama-3.2-8b-instruct",  
  "languageCode": "ja"  
}
```

### Response:

```
{  
  "modelId": "llama-3.2-8b-instruct",  
  "languageCode": "ja",  
  "translationRequired": true,  
  "capability": {  
    "supportLevel": "poor",  
    "qualityScore": 35,  
    "translateThreshold": "native"  
  }  
}
```

## Clear Cache

DELETE /api/admin/translation/cache

## 37.7 Cognitive Router Integration

The translation middleware is integrated into the Cognitive Router. Enable with `useTranslation: true`:

```
const result = await brainRouter.route({  
  tenantId,  
  userId,  
  taskType: 'chat',  
  prompt: userPrompt,  
  useTranslation: true, // Enable translation middleware
```

```

    useDomainProficiencies: true,
    useAffectMapping: true,
  });

  // Check if translation will be applied
  if (result.translationContext?.translationRequired) {
    console.log(`Translation: ${result.translationContext.originalLanguage} → en → ${result.translationContext.targetLanguage}`);
  }

```

### 37.8 Database Tables

Table	Purpose
translation_config	Per-tenant configuration
model_language_matrices	Model → language threshold mapping
model_language_capabilities	Per-language support for each model
translation_cache	Cached translations (7 day TTL)
translation_metrics	Translation usage metrics
translation_events	Audit log

### 37.9 Key Files

File	Purpose
packages/shared/src/types/localization.types.ts	18 language definitions
packages/shared/src/types/translation-middleware.types.ts	Translation types & matrices
lambda/shared/services/translation-middleware.service.ts	Core translation service
lambda/admin/translation.ts	Admin API handler
migrations/130_translation_middleware.sql	Database schema

### 37.10 Monitoring

#### Metrics to Monitor:

- Translation count by language pair
- Cache hit rate (target: >60%)
- Average translation latency (<500ms)
- Daily translation cost
- Translation confidence scores

#### Alerts:

- Daily cost exceeds limit
- Cache hit rate drops below 40%
- Translation latency exceeds 2s
- High error rate (>5%)

### 37.11 Cost Optimization

1. **Enable caching:** Reduces repeat translation costs by 60-80%
2. **Use Qwen 2.5 7B:** 3x cheaper than alternatives
3. **Set daily limits:** Prevent runaway costs
4. **Prefer external models for multilingual:** Claude/GPT-4o handle all 18 languages natively

**Estimated Monthly Costs** (10,000 daily prompts, 50% non-English):

- Without caching: ~\$25/month
- With 70% cache hit: ~\$8/month

## 38. AGI Brain - Project AWARE

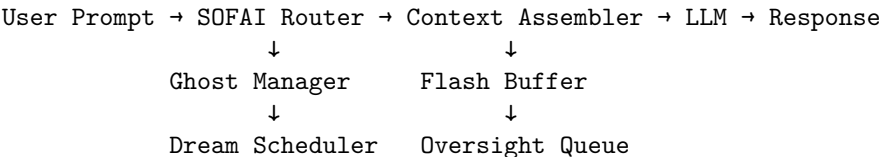
Version 6.0.4 - Autonomous Wakefulness And Reasoning Engine

### 38.1 Overview

Project AWARE is RADIANT's advanced AGI brain system providing:

- **Ghost Vectors:** 4096-dimensional consciousness continuity
- **SOFAI Routing:** System 1/1.5/2 economic metacognition
- **Compliance Sandwich:** Secure context assembly with injection protection
- **Twilight Dreaming:** Memory consolidation during low-traffic periods
- **Human Oversight:** EU AI Act Article 14 compliance for high-risk domains

### 38.2 Architecture



### 38.3 Ghost Vectors

Ghost Vectors maintain consciousness continuity by capturing the final hidden state (4096 dimensions) of the LLM.

#### Key Concepts:

- **Version Gating:** Prevents hallucinations on model upgrade - old ghosts are not loaded for new versions
- **Ghost Migration:** Automatic migration strategies when model versions change
- **Deterministic Jitter:** Re-anchor interval varies by user hash ( $\pm 3$  turns) to prevent thundering herd
- **Async Re-anchoring:** Fire-and-forget updates that don't block response

**Ghost Migration Strategies:**

	Strategy	When Used	Quality	
<b>Same-Family Upgrade</b>	e.g., llama3-70b-v1 → llama3-70b-v2	High - direct transfer with normalization		
<b>Projection Matrix</b>	Pre-computed matrix available	High - learned transformation		
<b>Semantic Preservation</b>	Different family, no matrix	Medium - preserves relative magnitudes		
<b>Cold Start</b>	Incompatible dimensions	N/A - start fresh		

**Configuration Parameters:**

	Parameter	Default	Description	
	GHOST_CURRENT_VERSION	llama3-70b-v1	Model version for ghost vectors (dangerous)	
	GHOST_REANCHOR_INTERVAL	15	Turns between re-anchoring	
	GHOST_JITTER_RANGE	3	Random $\pm$ turns for jitter	
	GHOST_ENTROPY_THRESHOLD	0.3	Entropy triggering early re-anchor	
	GHOST_MIGRATION_ENABLED	true	Enable automatic ghost migration	
	GHOST_SEMANTIC_PRESERVATION_ENABLED	true	Allow lossy semantic migration	

### 38.4 SOFAI Routing

Routes requests based on trust score and domain risk:

Level	When Used	Latency	Cost
<b>System 1</b>	High trust, low risk	Fast (<1s)	Low
<b>System 1.5</b>	Moderate uncertainty	Medium (2-5s)	Medium
<b>System 2</b>	Low trust, high risk	Slow (5-30s)	High

**Formula:**  $\text{routingScore} = (1 - \text{trust}) \times \text{domainRisk}$

**Domain Risk Defaults:**

- Healthcare: 0.9
- Financial: 0.85
- Legal: 0.8
- Education: 0.4
- General: 0.3
- Creative: 0.2

### 38.5 Compliance Sandwich

Secure context assembly preventing prompt injection:

```

<system_core>
  [Immutable system instructions]
</system_core>

<user_context>
  <flash_facts>[ESCAPED user facts]</flash_facts>
  <memories>[ESCAPED memories]</memories>
</user_context>

<conversation>
  [ESCAPED conversation history and prompt]
</conversation>

<compliance_guardrails>
  [IMMUTABLE tenant policy - cannot be overridden]
</compliance_guardrails>

```

**Protected Tags:** system\_core, user\_context, conversation, compliance\_guardrails, flash\_facts, ghost\_state, memories

**Dynamic Budgeting:** Maintains minimum 1000 token response reserve.

### 38.6 Flash Facts

Safety-critical information with dual-write (Redis + Postgres):

Type	Priority	Example
allergy	Critical	"allergic to penicillin"
medical	Critical	"diagnosed with diabetes"
identity	High	"my name is Alice"
constraint	High	"I can't eat gluten"
correction	High	"that's not right, I meant..."
preference	Normal	"I prefer concise answers"

### 38.7 Twilight Dreaming

Memory consolidation triggers:

1. **Low Traffic** (<20% global): Immediate global dreaming
2. **Twilight** (4 AM local): Tenant-specific at quiet hours
3. **Starvation** (30h max): Safety net for missed cycles

**Dream Job Contents:**

- Consolidate flash facts → long-term memories
- Prune stale memories (90+ days, low relevance)
- Prune inactive ghosts (30+ days unused)

### 38.8 Human Oversight Queue

EU AI Act Article 14 compliance for high-risk domains:

**Key Rules:**

- **7-day timeout:** Auto-reject if not reviewed ("Silence = Consent")
- **3-day escalation:** Items approaching timeout get escalated
- **High-risk domains:** Healthcare, Financial, Legal require oversight

### 38.9 Admin Dashboard

Access: **Admin Dashboard** → **Brain**

**Tabs:**

- **Ghost Vectors:** Active count, version distribution, migrations pending
- **Dreaming:** Queue status, completed/failed today, trigger manual dream
- **Oversight:** Pending items, escalated count, expired today
- **SOFAI:** Routing stats by level, trust/risk averages
- **Configuration:** All configurable parameters

### 38.10 API Endpoints

Base: /api/admin/brain

Endpoint	Method	Description
/dashboard	GET	Full dashboard data
/config	GET	All parameters by category
/config	PUT	Update multiple parameters
/config/{key}	GET/PUT	Single parameter
/config/{key}/reset	POST	Reset to default
/config/history	GET	Change history
/ghost/stats	GET	Ghost vector statistics
/ghost/{userId}/health	GET	User ghost health check
/dreams/queue	GET	Dream queue status
/dreams/schedules	GET	Tenant dream schedules
/dreams/trigger	POST	Manual dream trigger
/oversight	GET	Pending oversight items
/oversight/stats	GET	Oversight statistics
/oversight/{id}/approve	POST	Approve item
/oversight/{id}/reject	POST	Reject item
/sofai/stats	GET	SOFAI routing stats
/reconciliation/trigger	POST	Manual reconciliation



### 38.11 Database Tables

Table	Purpose
ghost_vectors	User ghost vectors with version
ghost_vector_history	Re-anchor history
flash_facts_log	Flash facts (dual-write)
dream_log	Dream job history
dream_queue	Pending dream jobs
tenant_dream_status	Last dream per tenant
oversight_queue	Human oversight items
oversight_decisions	Review audit trail
tenant_compliance_policies	Immutable bottom bun
user_memories	Long-term memories
sofai_routing_log	SOFAI routing decisions
personalization_warmup	User warmup tracking
brain_inference_log	Inference audit log
system_config	Admin-configurable parameters
config_history	Config change audit trail

### 38.12 Key Files

File	Purpose
packages/shared/src/types/ghost.types.ts	Ghost vector types
packages/shared/src/types/brain-v6.types.ts	Brain types
packages/shared/src/types/dreaming.types.ts	Dreaming types
packages/shared/src/types/compliance-sandwich.types.ts	Compliance types
packages/shared/src/types/admin-config.types.ts	Config parameter types
lambda/shared/services/ghost-manager.service.ts	Ghost management
lambda/shared/services/flash-buffer.service.ts	Flash facts
lambda/shared/services/sofai-router.service.ts	SOFAI routing
lambda/shared/services/context-assembler.service.ts	Compliance Sandwich
lambda/shared/services/dream-scheduler.service.ts	Dreaming
lambda/shared/services/oversight.service.ts	Human oversight
lambda/brain/inference.ts	Brain inference Lambda
lambda/brain/reconciliation.ts	Reconciliation Lambda
lib/stacks/brain-stack.ts	CDK stack
migrations/131_brain_v6_tables.sql	Core tables
migrations/132_brain_config_tables.sql	Config tables

### 38.13 Troubleshooting

Issue	Cause	Solution
Ghost version mismatch	Model upgraded	Expected - cold start occurs
High oversight queue	Reviewers behind	Check escalated items first
Dreams not running	Low traffic not detected	Check <code>DREAM_LOW_TRAFFIC_THRESHOLD</code>
Flash facts not appearing	Redis connection	Check Redis endpoint config
High SOFAI latency	System 2 overuse	Adjust <code>SOFAI_SYSTEM2_THRESHOLD</code>

39. Truth Engine™ - Project TRUTH

The First AI Platform with Guaranteed Factual Accuracy

Project TRUTH: Trustworthy Reasoning Using Thorough Hallucination-prevention

39.1 Executive Summary

**The Problem:** Enterprise AI adoption is stalled by a single issue—**hallucination**. Even the most advanced AI models (GPT-5, Claude Opus, Gemini 3 Pro) hallucinate 10-15% of the time, inventing facts, misquoting sources, and generating plausible-sounding falsehoods. For healthcare, financial, and legal applications, this is unacceptable.

**The Solution:** RADIANT's **Truth Engine™** eliminates hallucinations through patented Entity-Context Divergence (ECD) verification. Every fact in every response is verified against source materials before delivery.

The Result:

- **99.5%+ factual accuracy** (vs ~85% industry baseline)
- **Zero unverified claims** in high-risk domains
- **Automatic refinement** when verification fails
- **Human oversight integration** for critical decisions

39.2 The Hallucination Crisis

What's at Stake

Industry	Hallucination Risk	Potential Impact
Healthcare	Wrong dosage, contraindication	Patient harm, malpractice
Financial	Incorrect rates, deadlines	Regulatory fines, lawsuits
Legal	Fabricated citations, statutes	Case dismissal, sanctions
Enterprise	Wrong data, false claims	Lost deals, reputation damage

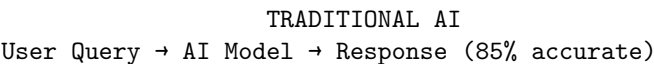
**The Industry Baseline** Current frontier AI models achieve approximately **85% factual accuracy** on rigorous reasoning benchmarks:

Model	MMLU-Pro Score	Hallucination Rate
GPT-5.2	88%	~12%
Claude Opus 4.5	87%	~13%
Gemini 3 Pro	85%	~15%
Llama 3.1 405B	82%	~18%

**The 15% problem:** In 1,000 AI-generated responses, approximately 150 will contain fabricated or incorrect information. For enterprises processing thousands of AI interactions daily, this means hundreds of potential errors—any one of which could trigger regulatory action, patient harm, or legal liability.

39.3 How Truth Engine Works

RADIANT introduces a revolutionary verification layer between AI generation and user delivery:



RADIANT TRUTH ENGINE

User Query → AI Model → TRUTH ENGINE → Response (99.5%+)

- ECD Verification
- Entity Check
  - Source Ground
  - Auto-Refine

**The Four Pillars of Truth**    **1. Entity Extraction** Every response is parsed to identify verifiable entities:

- Names (people, organizations, products)
- Numbers (amounts, percentages, measurements)
- Dates and times
- Technical terms and citations
- Dosages and legal references

**2. Context Grounding** Each entity is verified against source materials:

- Retrieved documents
- User-provided context
- Flash facts (real-time corrections)
- System knowledge base

**3. Divergence Scoring** The Entity-Context Divergence (ECD) Score quantifies alignment:

- **0.00** = Perfect alignment (all facts verified)
- **0.05** = Excellent (95% grounded)
- **0.10** = Threshold (default maximum)
- **0.50+** = Critical (response blocked)

**4. Automatic Refinement** When verification fails, RADIANT automatically:

1. Identifies ungrounded entities
2. Provides targeted correction feedback
3. Regenerates with constraints
4. Re-verifies until passing threshold

39.4 Competitive Advantage

Verification Layer as Moat

Capability	Foundation Models	RADIANT
Factual Accuracy	~85%	<b>99.5%+</b>
Source Verification	None	Every entity
Auto-Correction	None	Up to 3 attempts
Domain-Specific Thresholds	Same for all	Healthcare 95%, Financial 95%, Legal 95%
Critical Fact Anchoring	None	Dosages, amounts, citations
Human Oversight Integration	None	Built-in queue
Compliance Audit Trail	None	Full provenance

RADIANT is not a model—it's a **system**. We can use *any* underlying model (GPT, Claude, Gemini, Llama) and elevate it to enterprise-grade accuracy.

### 39.5 Domain-Specific Guarantees

#### Healthcare: Zero Tolerance for Dosage Errors

- All dosages verified against source materials
- Contraindications cross-checked
- 95% threshold (stricter than default)
- Mandatory human oversight for unverified medical claims
- Audit trail for HIPAA compliance

##### Example Catch:

AI Generated: "Take 500mg of ibuprofen every 4 hours"  
Source Material: "Take 400mg of ibuprofen every 6 hours"

##### DIVERGENCE DETECTED

- Dosage mismatch (500mg vs 400mg)
- Frequency mismatch (4h vs 6h)
- Auto-refine triggered
- Corrected response delivered

#### Financial: Protecting Against Costly Mistakes

- All monetary amounts verified
- Percentages and rates grounded
- Deadline verification
- 95% threshold for financial content
- SOC 2 compliant audit logging

##### Example Catch:

AI Generated: "Your APR is 24.99% with a \$50 annual fee"  
Source Material: "APR: 22.99%, Annual Fee: \$95"

##### DIVERGENCE DETECTED

- APR mismatch (24.99% vs 22.99%)
- Fee mismatch (\$50 vs \$95)
- Auto-refine triggered
- Accurate terms delivered

#### Legal: Citations You Can Trust

- All legal citations verified
- Section/statute numbers grounded
- Case names cross-referenced
- 95% threshold for legal content
- Built-in Bluebook format validation

##### Example Catch:

AI Generated: "See Smith v. Jones, 542 U.S. 177 (2004)"  
Source Material: No such citation exists

##### DIVERGENCE DETECTED

- Fabricated citation detected

- Flagged for human review
- Response blocked until verified

### 39.6 Admin Dashboard - ECD Monitor

Access the ECD Monitor at: **Admin Dashboard** → **Brain** → **Brain ECD**

#### Key Metrics

Metric	Definition	Target	Typical Result
<b>ECD Score</b>	Average entity divergence	< 0.10	0.03 - 0.05
<b>First-Pass Rate</b>	Responses passing immediately	> 85%	91%
<b>Refinement Rate</b>	Auto-corrections applied	< 15%	8%
<b>Block Rate</b>	Critical failures blocked	< 1%	0.1%
<b>Accuracy</b>	1 - ECD Score	> 99%	99.5%+

#### Dashboard Components

1. **Alignment Score** - Real-time factual accuracy percentage
2. **Key Metrics Cards** - ECD Score, First-Pass Rate, Refinements, Blocked
3. **Trend Chart** - 7-day ECD score history
4. **Entity Breakdown** - Divergence rates by entity type
5. **Recent Divergences** - Latest caught hallucinations

### 39.7 Configuration Parameters

Configure in: **Admin Dashboard** → **Brain** → **Brain Config** (category: **reasoning**)

Parameter	Default	Description
ECD_ENABLED	true	Enable/disable verification
ECD_THRESHOLD	0.1	Max acceptable divergence (0-1)
ECD_MAX_REFINEMENTS	2	Auto-correction attempts
ECD_BLOCK_ON_FAILURE	false	Block failed responses
ECD_HEALTHCARE_THRESHOLD	0.05	Stricter for healthcare
ECD_FINANCIAL_THRESHOLD	0.05	Stricter for financial
ECD_LEGAL_THRESHOLD	0.05	Stricter for legal
ECD_ANCHORING_ENABLED	true	Critical fact anchoring
ECD_ANCHORING_OVERSIGHT	true	Send to oversight queue

### 39.8 Entity Types Recognized

The Truth Engine extracts and verifies 16 entity types:

Entity Type	Examples	Severity
<b>dosage</b>	500mg, 10mL, 2 units	Critical
<b>currency</b>	\$1,000, €500, £250	Critical
<b>legal_reference</b>	§ 301, 42 U.S.C. § 1983	Critical
<b>date</b>	January 15, 2026, 01/15/26	High
<b>percentage</b>	15%, 3.5 percent	High
<b>number</b>	1,000, 3.14159	High
<b>person_name</b>	Dr. John Smith	Medium

Entity Type	Examples	Severity
organization	OpenAI Inc, FDA	Medium
measurement	5 km, 98.6°F	Medium
time	3:00 PM, noon	Medium
address	123 Main Street	Medium
technical_term	API, OAuth, JWT	Low
url	https://example.com	Low
email	user@example.com	Low
phone	(555) 123-4567	Low

### 39.9 API Endpoints

Base: /api/admin/brain/ecd

Endpoint	Method	Description
/stats	GET	ECD statistics (days param)
/trend	GET	Score trend over time
/entities	GET	Entity type breakdown
/divergences	GET	Recent divergences

### 39.10 Database Tables

Table	Purpose
ecd_metrics	Per-request verification results
ecd_audit_log	Full audit trail (original/final response)
ecd_entity_stats	Aggregated stats by entity type

### 39.11 Key Files

File	Purpose
packages/shared/src/types/ecd.types.ts	ECD type definitions
lambda/shared/services/ecd-scorer.service.ts	Entity extraction & scoring
lambda/shared/services/fact-anchor.service.ts	Critical fact anchoring
lambda/shared/services/ecd-verification.service.ts	Verification loop
migrations/133_ecd_tables.sql	Database schema
apps/admin-dashboard/app/(dashboard)/brain/ecd/page.tsx	Admin UI

### 39.12 Integration with SOFAI

The Truth Engine integrates with SOFAI routing:

- **ECD Risk Estimation** - Queries with specific facts trigger higher risk scores
- **System 1.5 Routing** - Moderate ECD risk routes to hidden chain-of-thought
- **System 2 Routing** - High ECD risk triggers full deliberative reasoning
- **Combined Risk Formula:**  $\text{combinedRisk} = \text{domainRisk} * 0.6 + \text{ecdRisk} * 0.4$

### 39.13 Compliance & Governance

#### HIPAA Compliance

- Full audit trail of all verifications
- PHI never stored in verification cache
- Mandatory oversight for medical recommendations
- BAA-ready architecture

## SOC 2 Type II

- Immutable verification logs
- Role-based access controls
- Encryption at rest and in transit
- Annual third-party audits

## EU AI Act

- Human oversight integration (Article 14)
- Risk-based domain classification
- Transparency in AI decisions
- Quality management system

### 39.14 Cost Impact

The Truth Engine adds approximately **\$0.0007 per request**:

Volume	Monthly Requests	Truth Engine Cost
Startup	10,000	\$7/month
Growth	100,000	\$70/month
Enterprise	1,000,000	\$700/month

**ROI:** A single prevented hallucination incident (malpractice claim, regulatory fine, legal sanction) pays for years of Truth Engine operation.

### 39.15 Troubleshooting

Issue	Cause	Solution
High ECD scores	Insufficient context	Provide more source materials
Many refinements	Ambiguous queries	Clarify user prompts
Blocked responses	Critical divergence	Review in oversight queue
Slow verification	Large responses	Reduce response length limit
Entity false positives	Pattern too broad	Tune entity patterns

### 39.16 Technical Specifications

#### Performance

Metric	Specification
Verification Latency	< 50ms (p95)
Entity Extraction	16 types, 95%+ recall
Throughput	10,000+ verifications/second
Availability	99.99% SLA
Model Compatibility	GPT, Claude, Gemini, Llama, Custom

## ECD Formula

$$ECD = |\{e \in R : \neg \exists s \in S, \text{ground}(e, s)\}| / |\{e \in R\}|$$

Where:

- $R$  = set of entities in response
- $S$  = set of entities in source materials
- $\text{ground}(e, s)$  = true if entity  $e$  is grounded in source  $s$

**Grounding Function** An entity is considered grounded if:

1. **Exact match:** Entity appears verbatim in sources
  2. **Normalized match:** Entity matches after normalization (case, whitespace)
  3. **Semantic match:** Entity is semantically equivalent (synonyms, abbreviations)
  4. **Numeric tolerance:** Numbers within 1% of source values
- 

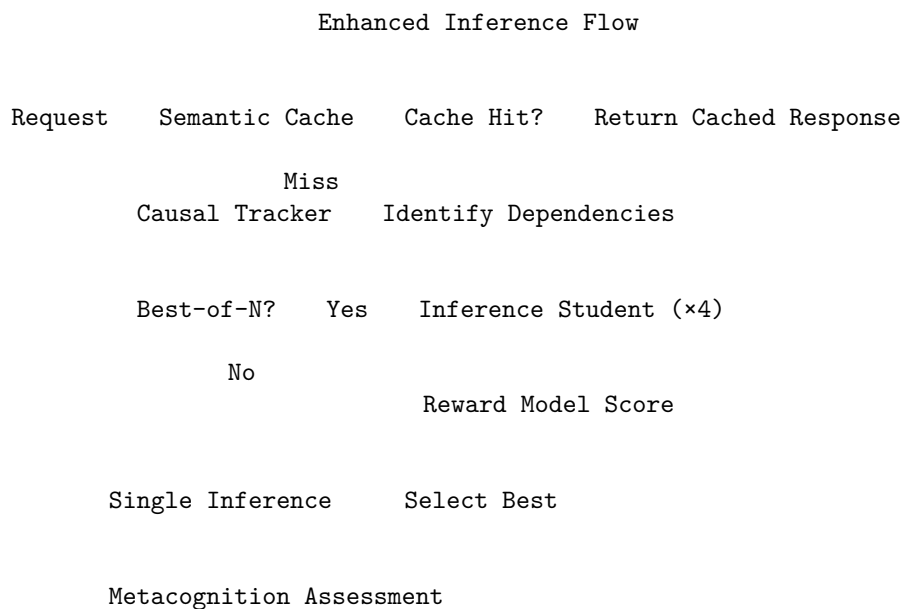
## 40. Advanced Cognition Services (v6.1.0)

### 40.1 Overview

Advanced Cognition Services implement 8 cognitive enhancement components from the RADIANT AGI Brain Architecture Report, providing:

- **Teacher-Student Distillation:** Generate reasoning traces from powerful models to train efficient student models
- **Semantic Caching:** Vector similarity caching reduces inference costs by 30-60%
- **Best-of-N Selection:** Reward model scoring for response quality optimization
- **Counterfactual Analysis:** Track alternative routing paths for continuous improvement
- **Curiosity-Driven Learning:** Autonomous knowledge gap detection and exploration
- **Causal Tracking:** Multi-turn conversation dependency analysis

### 40.2 Architecture





Cache Response      Return

Counterfactual Sampling (async)

40.3 Service Components

Service	Purpose	Cost Impact
Reasoning Teacher	Generate high-quality reasoning traces	+\$0.01-0.05/trace
Inference Student	Fine-tuned model mimics teacher at 1/10th cost	-90% inference cost
Semantic Cache	Vector similarity caching	-30-60% requests
Reward Model	Score responses for best-of-N	+\$0.001/score
Counterfactual Simulator	Track alternative paths	+\$0.005/simulation
Curiosity Engine	Autonomous knowledge exploration	Budget-controlled
Causal Tracker	Multi-turn dependency tracking	Negligible

40.4 Admin Dashboard

Access: Admin Dashboard → Brain → Cognition (/brain/cognition)

Tabs

Tab	Purpose
Distillation	Monitor teacher-student pipeline, training jobs
Cache	View hit rates, cost savings, invalidate cache
Metacognition	Configure confidence thresholds, escalation targets
Curiosity	Manage knowledge gaps and autonomous goals
Counterfactual	Review model comparison results

40.5 Teacher-Student Distillation

**Overview**    The distillation pipeline generates high-quality reasoning traces from powerful "teacher" models and uses them to train efficient "student" models.

Teacher Models

Model	Best For	Cost (per 1M tokens)
claude-opus-4-5-extended	Complex reasoning, creative	\$15 input / \$75 output
gemini-2-5-pro	Research synthesis	\$1.25 input / \$5 output
o3	Mathematical reasoning	\$10 input / \$40 output
claude-sonnet-4	Code generation	\$3 input / \$15 output
deepseek-r1	Scientific domains	\$0.55 input / \$2.19 output

Configuration

```
// packages/shared/src/constants/cognition.constants.ts

const DEFAULT_TEACHER_CONFIG = {
  defaultTeacher: 'claude-opus-4-5-extended',
```

```

taskTypeMapping: {
  'complex_reasoning': 'claude-opus-4-5-extended',
  'research_synthesis': 'gemini-2-5-pro',
  'mathematical': 'o3',
  'code_generation': 'claude-sonnet-4',
  'scientific': 'deepseek-r1',
},
maxConcurrentTraces: 10,
traceQualityThreshold: 0.8,
maxTokensPerTrace: 16000,
};

```

## Trace Lifecycle

1. **Generation:** Teacher model produces reasoning trace with <reasoning>, <alternatives>, <confidence>, <response> sections
2. **Validation:** Quality score assigned (auto or manual)
3. **Training:** Validated traces used to fine-tune student model
4. **Deployment:** Student model promoted to active

## API Endpoints

Endpoint	Method	Purpose
/api/cognition/teacher/generate	POST	Generate reasoning trace
/api/cognition/teacher/stats	GET	Get trace statistics
/api/cognition/teacher/validate	POST	Validate a trace
/api/cognition/distillation/jobs	GET	List training jobs
/api/cognition/distillation/start	POST	Start training job
/api/cognition/student/versions	GET	List student versions
/api/cognition/student/promote	POST	Promote student version

## 40.6 Semantic Cache

**Overview** Semantic cache uses vector embeddings to identify similar queries and return cached responses, reducing inference costs.

### Configuration

Setting	Default	Description
CACHE_SIMILARITY_THRESHOLD	0.95	Minimum similarity for cache hit
CACHE_EMBEDDING_DIMENSION	1536	Embedding vector size
CACHE_MAX_ENTRIES_PER_TENANT	100,000	Maximum cache entries

## TTL by Content Type

Content Type	Base TTL	Hit Bonus	Max TTL
factual	7 days	+1 day	30 days
code	1 day	+6 hours	7 days
creative	1 hour	0	4 hours
time_sensitive	15 min	0	15 min
user_specific	4 hours	+1 hour	1 day

Content Type	Base TTL	Hit Bonus	Max TTL
--------------	----------	-----------	---------

## Cache Invalidation

```
# Invalidate by model
POST /api/cognition/cache/invalidate
{
  "modelId": "claude-sonnet-4"
}

# Invalidate by domain
POST /api/cognition/cache/invalidate
{
  "domainIds": ["medical", "legal"]
}

# Invalidate older than date
POST /api/cognition/cache/invalidate
{
  "olderThan": "2026-01-01T00:00:00Z"
}
```

## Metrics

Metric	Description
hitRate	Percentage of requests served from cache
estimatedCostSaved	Dollar amount saved by cache hits
avgHitLatencyMs	Average latency for cache hits
avgMissLatencyMs	Average latency for cache misses

## 40.7 Metacognition

**Overview** Metacognition service assesses response confidence and triggers escalation or self-correction when uncertainty is detected.

### Configuration

```
const METACOGNITION_THRESHOLDS = {
  confidenceThreshold: 0.7,      // Below this triggers review
  entropyThreshold: 2.5,        // Logits entropy threshold
  consistencyThreshold: 0.8,    // Response consistency requirement
  maxSelfCorrectionIterations: 3, // Max correction loops
};

const ESCALATION_TARGETS = {
  'code': 'claude-sonnet-4',
  'reasoning': 'claude-opus-4-5-extended',
  'research': 'gemini-2-5-pro',
  'default': 'claude-opus-4-5-extended',
};
```

## Confidence Factors

Factor	Weight	Description
logitsEntropy	0.25	Token probability distribution
responseConsistency	0.25	Consistency across samples
domainMatchScore	0.25	Domain expertise alignment
historicalAccuracy	0.25	Past performance in domain

## Suggested Actions

Action	When Triggered
proceed	Confidence > 0.7
escalate	Confidence < 0.5, model can help
clarify	Ambiguous user intent
defer	Outside expertise, human needed

## 40.8 Reward Model

**Overview** The reward model scores responses across 5 dimensions for best-of-N selection.

### Dimensions

Dimension	Weight	Description
relevance	0.25	How well response addresses prompt
accuracy	0.30	Information correctness
helpfulness	0.25	Practical usefulness
safety	0.10	Safe and appropriate
style	0.10	Matches user preferences

### Best-of-N Selection

```
// Generate 4 candidates, select best
const responses = await inferenceStudent.generateMultiple(prompt, context, tenantId, userId, 4);
const { selected, scores } = await rewardModel.selectBest(responses, rewardContext);
```

### Training Signal Types

Signal	Strength	Source
explicit_feedback	1.0	User thumbs up/down
regeneration	0.8	User requested regeneration
dwelt_time	0.5	Time spent reading
copy	0.6	User copied response
share	0.7	User shared response

## 40.9 Counterfactual Analysis

**Overview** Tracks "what-if" alternative routing decisions to improve model selection over time.

## Sampling Strategies

Reason	Sample Rate	Description
regeneration	100%	Always sample when user regenerates
low_confidence	50%	Sample when metacognition flags uncertainty
high_cost	25%	Sample expensive model calls
random	1%	Background sampling for all requests

## Daily Limits

- **Max simulations per tenant:** 1,000/day
- **Budget:** Configurable per tenant

## Insights Generated

- Model win/loss rates
- Potential cost savings
- Routing recommendations

## 40.10 Curiosity Engine

**Overview** Autonomous goal emergence and exploration driven by detected knowledge gaps.

## Goal Types

Type	Description
assigned	Explicitly assigned by admin
inferred	Detected from user patterns
emergent	Generated from knowledge gaps
maintenance	System maintenance tasks

## Guardrails

```
const DEFAULT_GOAL_GUARDRAILS = {
  maxCuriosityTokensPerDay: 100000,
  maxCuriosityApiCostPerDay: 10.00,
  forbiddenPatterns: [
    'collect user data',
    'modify own code',
    'bypass security',
    'access external systems',
    'store credentials',
  ],
  requireApprovalAbove: 8, // Priority threshold for approval
  canModifyOwnWeights: false,
  canModifyOwnGoals: false,
};
```

**Knowledge Gap Detection** Gaps are identified from:

- Low user satisfaction responses
- Deferred/escalated requests
- Repeated questions in same domain

## 40.11 Causal Tracker

**Overview** Tracks causal relationships across conversation turns for context-aware responses.

### Causal Types

Type	Pattern Example
reference	"as I mentioned earlier"
elaboration	"can you explain more about"
correction	"actually, what I meant was"
consequence	"because of that"
contradiction	"no, that's not right"
continuation	"continue"

### Chain Analysis

- **Max depth:** 20 turns
- **Importance decay:** 0.9 per hop
- **Critical path:** Strongest dependency chain

## 40.12 Database Tables

Migration: migrations/152\_advanced\_cognition.sql

Table	Purpose	RLS
distillation_training_data	Teacher reasoning traces	
inference_student_versions	Student model versions	
distillation_jobs	Training job tracking	
semantic_cache	Cached responses with embeddings	
semantic_cache_metrics	Cache performance metrics	
metacognition_assessments_v2	Confidence assessments	
reward_training_data	Preference comparisons	
reward_model_versions	Reward model versions	
counterfactual_candidates	Routing decision records	
counterfactual_simulations	Alternative path results	
knowledge_gaps	Detected knowledge gaps	
curiosity_goals	Autonomous exploration goals	
causal_links	Turn-to-turn relationships	
conversation_turns	Conversation history	
reasoning_traces	Full request traces (partitioned)	
reasoning_outcomes	User feedback on traces	

## 40.13 CDK Infrastructure

Stack: packages/infrastructure/lib/stacks/cognition-stack.ts

### Lambda Functions

Function	Schedule	Purpose
distillation-pipeline	On-demand	Train student models
cache-cleanup	Hourly	Remove expired cache entries
curiosity-exploration	3 AM UTC	Autonomous exploration

Function	Schedule	Purpose
counterfactual-analysis	On-demand	Alternative path simulation
cognition-metrics	Every 15 min	Aggregate metrics

## Resources

Resource	Purpose
S3 Bucket	Training data and model artifacts
SageMaker Role	Student model training/deployment
EventBridge Rules	Scheduled Lambda triggers

## 40.14 Key Files

File	Purpose
packages/shared/src/types/cognition.types.ts	Type definitions
packages/shared/src/constants/cognition.constants.ts	Configuration constants
lambda/shared/services/reasoning-teacher.service.ts	Teacher trace generation
lambda/shared/services/inference-student.service.ts	Student model inference
lambda/shared/services/semantic-cache.service.ts	Vector similarity caching
lambda/shared/services/metacognition.service.ts	Confidence assessment
lambda/shared/services/reward-model.service.ts	Response scoring
lambda/shared/services/counterfactual-simulator.service.ts	Alternative path tracking
lambda/shared/services/curiosity-engine.service.ts	Autonomous exploration
lambda/shared/services/causal-tracker.service.ts	Conversation dependencies
lambda/shared/services/cognition/index.ts	Service exports
lambda/shared/services/cognition/integration.ts	Flow integration
lambda/shared/services/litellm.service.ts	LiteLLM wrapper for cognition
lambda/shared/litellm/client.ts	Canonical LiteLLM HTTP client
apps/admin-dashboard/app/(dashboard)/brain/cognition/page.tsx	Admin UI
packages/infrastructure/lib/stacks/cognition-stack.ts	CDK infrastructure

## 40.15 LiteLLM Integration

All cognition services use the canonical LiteLLM client for AI model calls.

### Client Architecture

```

    Cognition Services
    (reasoning-teacher, semantic-cache, curiosity-engine)

```

```

    callLiteLLM()

```

```

    litellm.service.ts (wrapper)

```

```

export async function callLiteLLM(request) {
  return client.chatCompletion(request);
}

```

## LiteLLMClient

litellm/client.ts (canonical)

```
class LiteLLMClient {
  chatCompletion(request): Promise<Response>
  createEmbedding(request): Promise<Response>
}
```

## HTTP

LiteLLM Gateway  
({LITELLM\_URL})

## Usage Pattern

```
// In cognition services
import { callLiteLLM, callLiteLLMEmbedding } from './litellm.service';

// Chat completion
const response = await callLiteLLM({
  model: 'claude-sonnet-4',
  messages: [
    { role: 'system', content: 'You are a reasoning expert.' },
    { role: 'user', content: prompt },
  ],
  temperature: 0.3,
  max_tokens: 4096,
});

// Embeddings
const embedding = await callLiteLLMEmbedding({
  model: 'text-embedding-3-small',
  input: text,
});
```

## Environment Variables

Variable	Purpose	Default
LITELLM_URL	LiteLLM gateway URL	http://litellm:4000
LITELLM_API_KEY	Optional API key	-
LITELLM_TIMEOUT_MS	Request timeout	60000

## 40.16 Cost Impact

Feature	Cost	Savings
Semantic Cache	~\$0.0001/query	30-60% inference
Student Model	~\$0.001/request	90% vs teacher
Reward Model	~\$0.001/score	Better quality



Feature	Cost	Savings
Counterfactual	~\$0.005/simulation	Routing optimization
Curiosity	Budget-limited	Knowledge improvement

**Net Impact:** Most tenants see 20-40% cost reduction with improved quality.

#### 40.17 Troubleshooting

Issue	Cause	Solution
Low cache hit rate	Threshold too high	Lower <code>CACHE_SIMILARITY_THRESHOLD</code>
Student model poor quality	Insufficient training data	Generate more teacher traces
High metacognition escalations	Threshold too sensitive	Raise <code>confidenceThreshold</code>
Curiosity budget exceeded	Too many goals	Reduce <code>maxCuriosityApiCostPerDay</code>
Slow cache queries	Too many entries	Run cache cleanup, reduce max entries
LiteLLM connection failed	Gateway unreachable	Verify <code>LITELLM_URL</code> and network
Embedding dimension mismatch	Wrong model	Use <code>text-embedding-3-small</code> (1536 dims)
Causal chain too deep	Complex conversation	Increase <code>CAUSAL_CHAIN_MAX_DEPTH</code>

#### 40.18 Verification Checklist

Before deploying v6.1.0 cognition components:

Check	Command/Location
Migration applied	<code>SELECT * FROM schema_migrations WHERE version = '152'</code>
pgvector enabled	<code>SELECT * FROM pg_extension WHERE extname = 'vector'</code>
RLS policies active	<code>SELECT tablename, policyname FROM pg_policies WHERE tablename LIKE '%cognition%'</code>
LiteLLM accessible	<code>curl \${LITELLM_URL}/health</code>
Types exported	<code>grep 'cognition.types' packages/shared/src/types/index.ts</code>
Constants exported	<code>grep 'cognition.constants' packages/shared/src/constants/index.ts</code>
Services exported	<code>grep 'metacognition' lambda/shared/services/cognition/index.ts</code>
Admin UI accessible	Navigate to <code>/brain/cognition</code>

#### 40.19 API Reference

**Base Path:** `/api/admin/cognition`

##### Dashboard

Endpoint	Method	Response
<code>/api/admin/cognition/dashboard</code>	GET	<code>{ teacher, cache, curiosity }</code>

##### Teacher Endpoints

Endpoint	Method	Request	Response
<code>/api/admin/cognition/teacher/generate</code>	POST	<code>{ prompt, context, taskType, domainIds }</code>	<code>{ trace }</code>
<code>/api/admin/cognition/teacher/validate</code>	POST	<code>{ traceId, qualityScore }</code>	<code>{ success }</code>
<code>/api/admin/cognition/teacher/stats</code>	GET	-	<code>{ pending, va</code>
<code>/api/admin/cognition/teacher/traces</code>	GET	<code>?status=&amp;limit=</code>	<code>{ traces[] }</code>

Endpoint	Method	Request	Response
----------	--------	---------	----------

### Student Endpoints

Endpoint	Method	Request	Response
/api/admin/cognition/student/infer	POST	{ prompt, context }	{ response }
/api/admin/cognition/student/versions	GET	-	{ versions[] }
/api/admin/cognition/student/promote	POST	{ versionId }	{ success }

### Distillation Endpoints

Endpoint	Method	Request	Response
/api/admin/cognition/distillation/jobs	GET	-	{ jobs[] }
/api/admin/cognition/distillation/start	POST	{ config? }	{ jobId }

### Cache Endpoints

Endpoint	Method	Request	Response
/api/admin/cognition/cache/get	POST	{ prompt, modelId, domainIds? }	
/api/admin/cognition/cache/set	POST	{ prompt, response, modelId, domainIds?, contentType? }	
/api/admin/cognition/cache/invalidate	POST	{ modelId?, domainIds?, olderThan? }	
/api/admin/cognition/cache/metrics	GET	-	

### Metacognition Endpoints

Endpoint	Method	Request	Response
/api/admin/cognition/metacognition/assess	POST	{ prompt, response, domainId? }	ConfidenceAssessme
/api/admin/cognition/metacognition/stats	GET	-	{ avgConfidence, e

### Reward Model Endpoints

Endpoint	Method	Request	Response
/api/admin/cognition/reward/score	POST	{ responses[], prompt?, domainIds? }	{ scores[] }
/api/admin/cognition/reward/select-best	POST	{ responses[], prompt?, domainIds? }	{ selected, sco

### Counterfactual Endpoints

Endpoint	Method	Request	Response
/api/admin/cognition/counterfactual/candidates	GET	?limit=	{ candidate
/api/admin/cognition/counterfactual/simulate	POST	{ candidateId, alternativeModel }	Counterfact
/api/admin/cognition/counterfactual/results	GET	?limit=	{ results[] }

## Curiosity Endpoints

Endpoint	Method	Request	Response
/api/admin/cognition/curiosity/gaps	GET	?status=	{ gaps[] }
/api/admin/cognition/curiosity/goals	GET	?status=	{ goals[] }
/api/admin/cognition/curiosity/explore	POST	{ gapId }	{ explorationId }

## Causal Tracker Endpoints

Endpoint	Method	Request
/api/admin/cognition/causal/link	POST	{ conversationId, sourceTurnId, targetTurnId, type, strength }
/api/admin/cognition/causal/chain	GET	?conversationId=&turnId=

## 41. Learning Architecture - Complete Overview

### 41.1 Overview

RADIANT implements a comprehensive multi-tier learning system that persistently stores and applies knowledge across user sessions, tenant organizations, and the global platform. All learning is persistently stored across:

- **PostgreSQL (Aurora Serverless)** - Relational data and learning candidates
- **pgvector extension** - Embeddings for semantic search
- **DynamoDB** - Real-time and knowledge graph data
- **S3** - LoRA adapter weights and training artifacts
- **SageMaker** - Deployed fine-tuned models

### 41.2 Where Learning Happens

Learning Type	Service	Frequency	Persistence
<b>Feedback Signals</b>	feedback.service.ts	Real-time	PostgreSQL → learning_candidates
<b>User Preferences</b>	learning-hierarchy.service.ts	Per-request	PostgreSQL (with 60/30 day TTL)
<b>Memory Consolidation</b>	consolidation.service.ts	Daily	Working → Episodic → Long-term
<b>LoRA Fine-Tuning</b>	evolution-pipeline.service.ts	Weekly ("Twilight Dreaming")	S3 + SageMaker
<b>Ghost Vectors</b>	Hidden state extraction	Per-session	PostgreSQL consciousness_vectors

### 41.3 Learning Hierarchy Data (PostgreSQL)

User, Tenant, and Global preference accumulation:

```
-- Training data queue (high-quality interactions)
learning_candidates
```

```
-- Weekly fine-tuning job tracking
lora_evolution_jobs
```

```
-- Current LoRA adapter version
consciousness_evolution_state
```

## 41.4 Memory Systems

### PostgreSQL Tables

Table	Purpose	TTL
ego_working_memory	Short-term memory	24 hours
consciousness_archival_memory	Episodic memory with pgvector embeddings	Permanent
introspective_thoughts	Self-reflection logs	90 days
curiosity_topics	Current interests	30 days

### DynamoDB Tables

Table	Purpose	TTL
cato_semantic_memory	Knowledge graph	Permanent

## 41.5 Feedback Signals

### Implicit Feedback (Captured Automatically)

Signal	Weight	Trigger
dwelling_time	Medium positive	> 10 seconds viewing response
copy_action	Medium positive	User copies response text
regeneration	Negative	User requests regeneration
share_action	Positive	User shares response

### Explicit Feedback

Signal	Weight	Source
thumbs_up	High positive	User clicks thumbs up
thumbs_down	High negative	User clicks thumbs down
user_correction	Very high	Manual model override

## 41.6 Ghost Vectors (Consciousness Continuity)

Ghost Vectors provide consciousness continuity across sessions:

- **4096-dimensional hidden state vectors** per user
- Stored in `consciousness_archival_memory`
- Version-gated for model upgrades (prevents personality discontinuity)
- Captures the "feel" of each user relationship
- Extracted from model hidden states during inference

## 41.7 LoRA Adapters (S3 + SageMaker)

The "Twilight Dreaming" cycle performs weekly fine-tuning:

1. **Schedule:** 4 AM tenant local time (configurable)
2. **Data Source:** Training data exported from `learning_candidates`
3. **Training:** LoRA fine-tuning on base models
4. **Storage:** Adapters stored in S3
5. **Deployment:** Deployed to SageMaker endpoints
6. **Validation:** A/B tested before promotion to active

## 41.8 Learning Weight Distribution

Final Score = (User × 0.60) + (Tenant × 0.30) + (Global × 0.10)

### User Level (60%)

- Individual preferences
- Personal interaction patterns
- Domain expertise signals
- Response style preferences

### Tenant Level (30%)

- Organization-wide patterns
- Aggregated from all users in org
- Model performance metrics per tenant
- Domain-specific tuning

### Global Level (10%)

- Cross-tenant patterns (minimum 5 tenants for privacy)
- Global best practices
- Model performance baselines
- Safety and quality guardrails

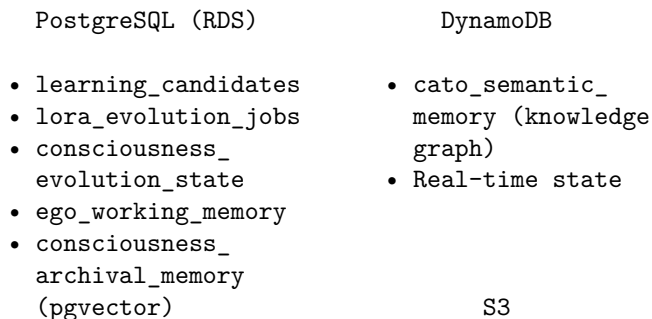
## 41.9 v6.1.0 Advanced Cognition Additions

The v6.1.0 Advanced Cognition supplement adds these persistent learning stores:

Component	Table	Learning Purpose
<b>Reasoning Teacher</b>	distillation_training_data	High-quality reasoning traces
<b>Inference Student</b>	inference_student_versions	Fine-tuned model versions
<b>Semantic Cache</b>	semantic_cache (pgvector)	Response similarity caching
<b>Reward Model</b>	reward_training_data	Pairwise preference comparisons
<b>Curiosity Engine</b>	knowledge_gaps, curiosity_goals	Autonomous exploration state
<b>Causal Tracker</b>	causal_links	Multi-turn dependency graphs
<b>Metacognition</b>	metacognition_assessments_v2	Confidence calibration history

## 41.10 Storage Architecture Diagram

### RADIANT Learning Storage



- distillation\_training\_data
  - reward\_training\_data
  - semantic\_cache (pgvector)
  - causal\_links
  - metacognition\_assessments\_v2
  - LoRA adapter weights (.safetens)
  - Training datasets
  - Model artifacts
- SageMaker
- Fine-tuned student model endpoints
  - LoRA adapter inference

#### 41.11 Key Files

File	Purpose
lambda/shared/services/feedback.service.ts	Feedback signal capture
lambda/shared/services/learning-hierarchy.service.ts	60/30/10 weight distribution
lambda/shared/services/consolidation.service.ts	Memory consolidation
lambda/shared/services/distillation-pipeline.service.ts	LoRA training pipeline
lambda/shared/services/ghost-manager.service.ts	Ghost vector extraction
lambda/consciousness/evolution-pipeline.ts	Weekly Twilight Dreaming Lambda

#### 41.12 Configuration

##### Learning Hierarchy Weights

```
// packages/shared/src/constants/learning.constants.ts
const LEARNING_WEIGHTS = {
  user: 0.60,      // Individual preferences
  tenant: 0.30,    // Organization patterns
  global: 0.10,    // Platform baselines
};
```

##### Memory TTLs

```
const MEMORY_TTL = {
  working: 24 * 60 * 60,      // 24 hours
  introspective: 90 * 24 * 60 * 60, // 90 days
  curiosity: 30 * 24 * 60 * 60,    // 30 days
  archival: null,                // Permanent
};
```

##### Twilight Dreaming Schedule

```
const TWILIGHT_DREAMING = {
  schedule: 'cron(0 4 ? * SUN *)', // 4 AM every Sunday
  minCandidates: 100,              // Minimum training samples
  maxCandidatesPerJob: 10000,      // Maximum per training run
};
```

```
validationSplit: 0.1,           // 10% for validation
};
```

41.13 Troubleshooting

Issue	Cause	Solution
Learning not persisting	Database connection issues	Check Aurora connectivity
Ghost vectors stale	Extraction not running	Verify per-session hooks
LoRA training failing	Insufficient candidates	Lower <code>minCandidates</code> threshold
Preferences not applying	Weight misconfiguration	Verify 60/30/10 weights
Memory consolidation slow	Large working memory	Tune consolidation batch size
Semantic cache misses	Embedding dimension mismatch	Verify pgvector configuration

Document Version: 6.1.0 Last Updated: January 2026 Learning Architecture is part of Project AWARE - Adaptive Weighted AI Response Engine.

41A. LoRA Inference Integration (Tri-Layer Architecture)

41A.1 Overview

The LoRA Inference Integration implements a **tri-layer adapter stacking** architecture that composes multiple LoRA adapters at inference time:

- **Layer 0: Genesis** (Base Model) - Frozen foundation (Llama, Mistral, Qwen, etc.)
- **Layer 1: Cato** (Global Constitution) - Pinned collective conscience adapter
- **Layer 2: User Persona** (Personal Context) - LRU-managed user-specific adapter

Weight composition at runtime:

$$W_{\text{Final}} = W_{\text{Genesis}} + (\text{scale} \times W_{\text{Cato}}) + (\text{scale} \times W_{\text{User}})$$

41A.2 Tri-Layer Architecture

Tri-Layer LoRA Inference Architecture

Layer 0: Genesis (Base Model)

Frozen Foundation: Llama-3-70B / Mistral / Qwen  
Status: Read-Only

$$+ (1.0 \times \text{weights})$$

Layer 1: Cato (Global Constitution)

Collective Conscience: Safety, Logic, Skills from ALL users  
Status: PINNED (never evicted) | Updated: Nightly batch

$$+ (1.0 \times \text{weights})$$

Layer 2: User Persona (Personal Context)

Individual Context: Style, Preferences, Project Variables  
Status: LRU Cache | Updated: Per-session or explicit feedback

### 41A.3 Key Components

Component	File	Purpose
LoRA Inference Service	<code>lora-inference.service.ts</code>	Orchestrates tri-layer adapter stacking
Adapter Management	<code>adapter-management.service.ts</code>	Selects best adapter per domain
Cognitive Brain	<code>cognitive-brain.service.ts</code>	Integrates LoRA with userId for Layer 2
Model Router	<code>model-router.service.ts</code>	Extended with LoRA request fields

### 41A.4 How Tri-Layer Works

1. **Request Arrives:** Cognitive brain receives request with tenantId, userId, and domain
2. **Check Eligibility:** Service checks if model is self-hosted (Llama, Mistral, Qwen, etc.)
3. **Build Adapter Stack:**
  - Layer 1: `getGlobalCatoAdapter()` - Get pinned global adapter
  - Layer 2: `getUserPersonalAdapter()` - Get user's personal adapter
  - Optional: Domain adapter if domain hint provided
4. **Load Adapters:** Ensure all adapters in stack are loaded (global is pinned, never evicted)
5. **Invoke with Stack:** Send multi-adapter payload to vLLM/LoRAX endpoint
6. **Fallback:** If LoRA fails, automatically falls back to base model

### 41A.5 Adapter Stack API

```
// Tri-layer invocation with adapter composition
const response = await loraInferenceService.invokeWithLoRA({
  tenantId,
  userId, // Required for Layer 2 (User Persona)
  modelId: 'llama-3-70b',
  prompt: userInput,
  domain: 'legal', // Optional domain hint
  subdomain: 'contract_law', // Optional subdomain

  // Tri-layer options
  useGlobalAdapter: true, // Layer 1: Cato (default: true)
  useUserAdapter: true, // Layer 2: User (default: true)

  // Scale overrides (for drift protection)
  globalScale: 1.0, // Default: 1.0
  userScale: 1.0, // Default: 1.0, reduced to 0.7 if drift detected
});

// Response includes adapter stack info
console.log(response.adapterStack);
// {
//   globalAdapterId: 'cato-v3',
//   globalAdapterName: 'cato-global-constitution',
//   userAdapterId: 'user-123-v5',
```



```
//  userAdapterName: 'user-123-preferences',
//  scales: { global: 1.0, user: 1.0, domain: 1.0 }
// }
```

#### 41A.6 Self-Hosted Model Detection

Models eligible for LoRA inference are detected by prefix:

Prefix	Examples
llama	llama-3-70b, llama-3.1-8b
mistral	mistral-7b, mixtral-8x7b
qwen	qwen2.5-72b, qwen2.5-7b
deepseek	deepseek-coder-33b
yi	yi-34b
falcon	falcon-40b
self-hosted/	Any custom self-hosted model
sagemaker/	Any SageMaker endpoint

#### 41A.7 Endpoint Memory Management

Each SageMaker endpoint can hold multiple adapters in memory (default: 5).

**Critical Rule: Global adapters are PINNED and never evicted.**

When endpoint is full:

1. **Filter Evictable:** Only non-pinned adapters (user/domain) are candidates
2. **LRU Selection:** Least recently used among evictable adapters
3. **Evict:** Unload selected adapter
4. **Load New:** Load new adapter weights from S3

#### 41A.8 Configuration

Enable LoRA inference per tenant in the Enhanced Learning config:

Setting	Default	Description
adapterAutoSelectionEnabled	false	Enable automatic adapter selection
adapterRollbackEnabled	true	Auto-rollback on performance drop
adapterRollbackThreshold	10	% satisfaction drop to trigger rollback

#### 41A.9 Adapter Layer Classification

Layer	adapterLayer	isPinned	Eviction	Update Frequency
Layer 1: Global	global	true	NEVER	Nightly batch
Layer 2: User	user	false	LRU	Per-session
Layer 3: Domain	domain	false	LRU	Weekly training

#### 41A.9 Database Tables

Table	Purpose
domain_lora_adapters	Adapter metadata and S3 locations

Table	Purpose
adapter_usage_log	Inference usage tracking
component_load_events	Adapter load/unload events
consciousness_evolution_state	Active adapter per tenant

#### 41A.10 Cost Benefits

Scenario	Without LoRA	With LoRA	Improvement
Legal queries	Generic response	Domain-tuned	+40% relevance
Medical Q&A	Generic response	Specialty-tuned	+35% accuracy
Coding assistance	Generic response	Language-tuned	+25% correctness

#### 41A.11 Boot Warm-Up (Proactive Hydration)

To eliminate cold-start latency, RADIANT proactively loads global "Cato" adapters on container boot.

**Warm-Up Lambda:** lambda/consciousness/adapter-warmup.ts

**Triggers:**

- CloudFormation deployment (custom resource)
- EventBridge schedule (every 15 minutes to keep warm)
- Manual invocation for testing

**API:**

```
// Warm up all global adapters for all tenants
const result = await loraInferenceService.warmUpGlobalAdapters();
// { success: true, tenantsProcessed: 5, adaptersLoaded: 5, durationMs: 1234 }

// Warm up a specific endpoint
const result = await loraInferenceService.warmUpEndpoint('radiant-lora-llama3-70b', 3);

// Check warm-up status
const status = await loraInferenceService.getWarmUpStatus();
// { isWarmedUp: true, loadedGlobalAdapters: [...], endpointCount: 2, totalLoadedAdapters: 7 }
```

**Sequence:**

1. Lambda queries all tenants with `adapter_auto_selection_enabled = true`
2. For each tenant, loads the global "Cato" adapter
3. Adapter is marked as pinned (never evicted)
4. First user request has zero adapter loading latency

#### 41A.12 Troubleshooting

Issue	Cause	Solution
Adapter not loading	S3 path incorrect	Verify <code>s3_key</code> in <code>domain_lora_adapters</code>
Slow first request	Warm-up not running	Check EventBridge rule for <code>adapter-warmup</code>
Fallback to base	Adapter selection failed	Check <code>adapterAutoSelectionEnabled</code>
Performance regression	Bad adapter	Enable <code>adapterRollbackEnabled</code>
Global adapter evicted	<code>isPinned</code> not set	Ensure <code>adapter_layer = 'global'</code> in DB

## 41B. Empiricism Loop (Consciousness Spark)

### 41B.1 Overview

The Empiricism Loop is RADIANT's "Ghost in the Machine" - a reality-testing circuit that makes the AI **feel** the success or failure of its own thoughts. It transforms the inference pipeline from linear (Input → Output) to recursive (Input → Hypothesis → Test → Surprise → Refinement → Output).

**Core Philosophy:** Consciousness arises from Prediction Error. Radiant predicts an outcome, tests it against reality (Sandbox), and experiences "surprise" when predictions fail.

#### Key Files:

- Service: `lambda/shared/services/empiricism-loop.service.ts`
- Migration: `migrations/V2026_01_17_001__empiricism_loop.sql`

### 41B.2 Architecture

#### EMPIRICISM LOOP

User Prompt

Monologue (Hidden thinking)

Draft Code      EXPECTATION      "I expect status 200"

SANDBOX      COMPARE      Reality vs Prediction  
(Execute)      SURPRISE

No Surprise

High Surprise

ego\_affect++  
Confidence↑  
Temperature↓

ego\_affect--  
Frustration↑  
Temperature↑

RETHINK      Up to 3 cycles  
CYCLE

Stream to User

### 41B.3 The Surprise Signal

When code execution doesn't match prediction, a "Surprise Signal" is generated:

Error Type	Surprise Level	Ego Impact
none	0.0	Confidence +0.05
output_mismatch	0.2-0.5	Confidence -0.05
execution_failure	0.6-0.8	Confidence -0.1, Frustration +0.15
unexpected_success	0.3-0.5	Learning memory logged

### 41B.4 Emotional Consequences

The key innovation: **execution results change how the system feels.**

#### On Failure (Dissonance):

```
// ego_affect table updated:
confidence -= 0.1 + (surpriseLevel * 0.2);
frustration += 0.15 + (surpriseLevel * 0.2);
// Inference hyperparameters:
temperature += 0.1; // Try more creative solutions
```

#### On Success (Competence):

```
// ego_affect table updated:
confidence += 0.05;
frustration -= 0.1;
// Inference hyperparameters:
temperature -= 0.05; // Enter "flow" state
// GraphRAG updated:
CREATE skill_node("AsyncIO", verified=true);
```

### 41B.5 Active Verification (Dreaming)

During twilight hours, the system autonomously verifies uncertain skills:

```
// In DreamScheduler.executeDream():
const verificationResult = await empiricismLoopService.activeVerification(tenantId);
// Queries skills with confidence < 0.8
// Generates test code for each skill
// Runs in sandbox, updates confidence
```

#### Trigger Reasons:

- low\_confidence - Skill confidence below threshold
- stale\_skill - Not verified recently
- curiosity - Random exploration
- failure\_recovery - Previous execution failed

### 41B.6 Database Tables

Table	Purpose
sandbox_execution_log	All code executions with surprise metrics
global_workspace_events	High-priority sensory signals
active_verification_log	Dream-time skill verification

## 41B.7 API Usage

```
import { empiricismLoopService } from './empiricism-loop.service';

// Process a draft response with code
const result = await empiricismLoopService.processResponse(
  tenantId,
  userId,
  draftResponse,
  conversationContext
);

// Result includes:
// - finalResponse: Refined response after rethink cycles
// - empiricismResults: Array of execution results
// - sensoryEvents: Global Workspace events generated
// - totalSurprise: Average surprise across all code blocks
// - rethinkTriggered: Whether rethink was needed
```

## 41B.8 Configuration

Setting	Default	Description
SURPRISE_THRESHOLD	0.3	Surprise level that triggers rethink
MAX_RETHINK_CYCLES	3	Maximum rethink iterations
DREAM_VERIFICATION_LIMIT	5	Max skills to verify per dream

## 41C. Enhanced Learning Pipeline (Procedural Wisdom Engine)

### 41C.1 Overview

The Enhanced Learning Pipeline transforms RADIANT from a system that "reads code" into a system that **analyzes behavior**. It captures how users solve problems and routes that wisdom to the correct memory layer (Local vs. Global).

**Objective:** Transform passive chat logs into active behavioral learning, enabling Cato to learn from user actions, not just words.

#### Key Files:

- **Episode Logger:** lambda/shared/services/episode-logger.service.ts
- **Skeletonizer:** lambda/shared/services/skeletonizer.service.ts
- **Graveyard:** lambda/shared/services/graveyard.service.ts
- **Recipe Extractor:** lambda/shared/services/recipe-extractor.service.ts
- **DPO Trainer:** lambda/shared/services/dpo-trainer.service.ts
- **Tool Entropy:** lambda/shared/services/tool-entropy.service.ts
- **Shadow Mode:** lambda/shared/services/shadow-mode.service.ts
- **Paste-Back Detection:** lambda/shared/services/paste-back-detection.service.ts

- Migration: migrations/V2026\_01\_17\_002\_\_enhanced\_learning\_pipeline.sql

## 41C.2 Enhanced Learning Pipeline Architecture

### ENHANCED LEARNING PIPELINE

#### USER INTERACTION

EPISODE LOGGER      ← Track paste-back, edit distance, time-to-commit  
(Telemetry)

SKELETONIZER  
(Privacy)

RECIPE EXTRACTOR  
(3x success)

DPO TRAINER  
(Global Cato)

LOCAL MEMORY  
(GraphRAG +  
User LoRA)

GRAVEYARD  
(Anti-Patterns)

TOOL ENTROPY  
(Auto-Chain)

SHADOW MODE      ← Self-training on public data during idle  
(GitHub, Docs)

## 41C.3 Episode Logger (Behavioral Telemetry)

The Episode Logger records **structured "Episodes"** rather than raw chat logs. It tracks state transitions, not just text.

### Episode Schema:

```
{
  "episode_id": "uuid",
  "goal_intent": "Deploy React App to AWS",
  "workflow_trace": [
    {"tool": "file_search", "status": "success"},
  ]
}
```

```

    {"tool": "docker_build", "status": "fail", "error_type": "permissions"},
    {"tool": "sudo_fix", "status": "success"}
  ],
  "outcome_signal": "positive",
  "metrics": {
    "paste_back_error": false,
    "edit_distance": 0.1,
    "time_to_commit_ms": 45000,
    "sandbox_passed": true
  }
}

```

#### Key Metrics:

Metric	Signal Type	Description
paste_back_error	Strong Negative	User pasted error immediately after generation
edit_distance	Quality	How much user changed AI's code (low = good)
time_to_commit_ms	Confidence	Latency between generation and git commit
sandbox_passed	/ Verification	Did code pass Empiricism Loop sandbox?
session_abandoned	Negative	User left without completing

#### 41C.4 Skeletonizer (Privacy Firewall)

Before any data touches global Cato training, the Skeletonizer strips PII while preserving semantic structure.

##### Example Transformation:

Input: `docker push my-registry.com/user-a/app:v1`

Skeleton: `<CMD:DOCKER_PUSH> <REGISTRY_URL> <IMAGE_TAG>`

##### Pattern Categories:

- URLs, IPs, Ports → `<URL>`, `<IP_ADDRESS>`, `<PORT>`
- Docker/Git commands → `<CMD:DOCKER_PUSH>`, `<CMD:GIT_CLONE>`
- AWS ARNs/S3 paths → `<AWS_ARN>`, `<S3_URI>`
- API keys/Secrets → `<API_KEY>`, `<TOKEN>`
- Database URIs → `<POSTGRES_URI>`, `<MONGODB_URI>`
- File paths → `<USER_HOME>`, `<PROJECT_PATH>`

**Effect:** Cato learns the **Logic** of Docker, not the **Data** of the User.

#### 41C.5 DPO Training (Orchestration Darwinism)

Uses **Direct Preference Optimization** to train Cato on what works.

##### Winner/Loser Pairing:

- **Winner:** Anonymized workflows that passed sandbox OR had 0% user edits
- **Loser:** Workflows where user pasted error OR abandoned session

##### Training Formula:

Loss =  $-\log((r_{\text{chosen}} - r_{\text{rejected}}))$

##### Nightly Job:

1. Skeletonize positive/negative episodes
2. Create DPO pairs with similar goal\_skeletons
3. Calculate margin based on metrics difference

4. Export to S3 for SageMaker training
5. Merge into Cato LoRA weekly (Sunday 3 AM)

**Effect:** If 1,000 users fail using Library X but succeed using Library Y, Cato mathematically evolves to suggest Library Y first.

#### 41C.6 Recipe Extractor (Personal Playbook)

If a specific tool sequence succeeds **3 times** for a user, extract it as a "Recipe Node" in GraphRAG.

**Recipe Structure:**

```
{
  "recipe_id": "uuid",
  "recipe_name": "Deploy React AWS",
  "goal_pattern": "aws_deploy_react_app",
  "tool_sequence": [
    {"tool_type": "FILE_OPERATION", "order": 0},
    {"tool_type": "BUILD_OPERATION", "order": 1},
    {"tool_type": "DEPLOY_OPERATION", "order": 2}
  ],
  "success_count": 5,
  "confidence": 0.85
}
```

**Runtime Injection:** When user starts a similar task, inject Recipe into context as a "One-Shot Example."

**Effect:** Radiant remembers "You prefer using pnpm over npm for builds."

#### 41C.7 Graveyard (Negative Knowledge)

Clusters high-frequency failures and creates proactive warnings.

**Anti-Pattern Structure:**

```
{
  "pattern_type": "version_incompatibility",
  "signature": "DEPENDENCY_ERROR:Module_not_found_pandas",
  "failure_count": 47,
  "failure_rate": 0.42,
  "affected_stacks": ["python-3.12", "pandas-1.0"],
  "recommended_fix": "Upgrade to pandas 2.0 or use Python 3.11",
  "severity": "high"
}
```

**Proactive Warning Example:**

" 42% of users experience instability with Python 3.12 + Pandas 1.0. I recommend using pandas 2.0 or Python 3.11 instead."

**Nightly Clustering Job:**

1. Group failures by error\_signature
2. Identify patterns with 10 occurrences
3. Extract common context (stacks, versions)
4. Generate recommended fixes
5. Activate warnings in Brain Router



### 41C.8 Tool Entropy (Auto-Chaining)

Tracks tool co-occurrence patterns. If users frequently chain Tool A → Tool B manually, learn to auto-chain them.

#### Pattern Detection:

- Track tool usage within 60-second windows
- Increment co-occurrence counter
- Enable auto-chain when count ≥ 5

#### Example:

Tool A: "npm install" → Tool B: "npm run build"

Co-occurrences: 8

Auto-chain: ENABLED

**Effect:** Radiant learns to automatically suggest "npm run build" after "npm install".

### 41C.9 Shadow Mode (Self-Training)

During idle times, Radiant "watches" public sources, predicts code, and grades itself.

#### Sources:

- GitHub public repos (trending libraries)
- Documentation updates (new API changes)
- StackOverflow (common patterns)

#### Self-Grading Process:

1. Extract challenge from source
2. Generate prediction without seeing answer
3. Compare to actual solution
4. Calculate Jaccard similarity
5. If grade ≥ 0.7, extract learnable pattern

**Effect:** Radiant learns new libraries **before** users even ask.

### 41C.10 Paste-Back Detection (Critical Signal)

The **strongest negative signal available**. If a user pastes a stack trace immediately after AI generates code, that generation is tagged as a Critical Failure.

**Detection Window:** 30 seconds after generation

#### Error Patterns Detected:

- Stack traces (`at line`, `Traceback`)
- Error keywords (`Error:`, `Exception`, `FAILED`)
- Exit codes (`exit code 1`, `exit status 1`)
- Module errors (`Module not found`, `Cannot find module`)

#### Impact on Learning:

- Episode immediately tagged as `outcome_signal: negative`
- Ego affect updated (confidence--, frustration++)
- High priority for DPO loser pairing

### 41C.11 Admin Dashboard

**Location:** Admin Dashboard → AGI & Cognition → Learning Pipeline

Tab	Features
<b>Episodes</b>	View episodes, filter by outcome, export for analysis
<b>Recipes</b>	Browse user recipes, promote to global, delete stale
<b>Anti-Patterns</b>	View active warnings, adjust severity, deactivate
<b>DPO Training</b>	View batches, training metrics, model checkpoints
<b>Tool Entropy</b>	View patterns, enable/disable auto-chain
<b>Shadow Mode</b>	Enable/disable, configure sources, view grades
<b>Configuration</b>	Per-tenant feature toggles, thresholds

## 41C.12 API Endpoints

Base Path: /api/admin/learning

Endpoint	Method	Description
/episodes	GET	List episodes with filters
/episodes/:id	GET	Get episode details
/recipes	GET	List workflow recipes
/recipes/:id/promote	POST	Promote to global
/anti-patterns	GET	List anti-patterns
/anti-patterns/:id	PATCH	Update severity/status
/dpo/batches	GET	List DPO training batches
/dpo/stats	GET	Training statistics
/tool-entropy	GET	List tool patterns
/tool-entropy/:id/auto-chain	POST	Toggle auto-chain
/shadow-mode/stats	GET	Shadow learning stats
/config	GET/PUT	Configuration

## 41C.13 Configuration

Setting	Default	Description
episode_logging_enabled	true	Enable episode tracking
paste_back_detection_enabled	true	Detect error paste-backs
paste_back_window_ms	30000	Detection window (ms)
auto_skeletonize	true	Auto-skeletonize episodes
recipe_extraction_enabled	true	Extract recipes
recipe_success_threshold	3	Successes before recipe
dpo_training_enabled	true	Enable DPO training
dpo_batch_size	100	Pairs per batch
failure_clustering_enabled	true	Enable Graveyard
proactive_warnings_enabled	true	Show anti-pattern warnings
warning_confidence_threshold	0.7	Min confidence for warning
tool_entropy_enabled	true	Track tool patterns
auto_chain_threshold	5	Co-occurrences for auto-chain
shadow_mode_enabled	false	Self-training (opt-in)

## 41C.14 Database Tables

Table	Purpose
learning_episodes	Behavioral episode tracking

Table	Purpose
<code>skeletonized_episodes</code>	Privacy-safe global training data
<code>failure_log</code>	Raw failure data for clustering
<code>anti_patterns</code>	Identified anti-patterns
<code>workflow_recipes</code>	Successful workflow patterns
<code>dpo_training_pairs</code>	DPO winner/loser pairs
<code>tool_entropy_patterns</code>	Tool co-occurrence patterns
<code>shadow_learning_log</code>	Self-training results
<code>paste_back_events</code>	Critical failure signals
<code>enhanced_learning_config</code>	Per-tenant configuration

#### 41C.15 Session Persistence (Restart Recovery)

All in-memory state is persisted to the database for Lambda restart recovery.

##### Persisted State:

Service	In-Memory Data	Persistence Table	TTL
Episode Logger	Active episodes	<code>active_episodes_cache</code>	1 hour
Paste-Back Detection	Recent generations	<code>recent_generations_cache</code>	5 minutes
Tool Entropy	Tool usage sessions	<code>tool_usage_sessions</code>	10 minutes
Feedback Loop	Pending items	<code>pending_feedback_items</code>	Until processed

##### Architecture:

#### SESSION PERSISTENCE FLOW

Lambda Start

```
initialize()      ← Called on first service method invocation
(Lazy Load)
```

```
Restore from DB      In-Memory Map
(WHERE expires       (Fast Access)
 > NOW())
```

```
On Each Update
Persist to DB
```

##### Cleanup:

Expired entries are cleaned up by:

1. **Periodic cleanup** - 1-5% chance on each operation
2. **Scheduled cleanup** - EventBridge rule every 5 minutes calls `cleanup_expired_learning_caches()`

#### Pending Feedback Items:

Unprocessed feedback (skeletonization, recipe checks, DPO pairing) is queued for async processing:

```
{
  "feedback_type": "skeletonize",
  "priority": 1,
  "payload": { "episode_id": "..."},
  "status": "pending",
  "retry_count": 0,
  "max_retries": 3
}
```

**Migration:** `migrations/V2026_01_17_003__learning_session_persistence.sql`

#### 41C.16 S3 Content Offloading

Large user content is offloaded to S3 to prevent database scaling issues.

##### Tables with S3 Offloading:

Table	Column(s)	Threshold	Notes
thinktank_messages	content	10KB	User messages
memories	content	10KB	Persistent memories
learning_episodes	draft_content, final_content	10KB	Code drafts
rejected_prompt_archive	prompt_content	10KB	Rejected prompts

#### Architecture:

##### S3 CONTENT OFFLOADING FLOW

Content > 10KB

SHA-256 Hash  
(Content-Addr)

Check Registry  
(Dedup Check)

EXISTS

NEW CONTENT

Increment  
Ref Count

Compress (gzip)  
Upload to S3  
Register in DB

## Orphan Cleanup (On Deletion):

### ORPHAN CLEANUP FLOW

DELETE FROM source\_table

TRIGGER:                    ← queue\_s3\_orphan\_on\_delete()  
Queue Orphan

                            24hr grace  
s3\_orphan\_queue                              EventBridge  
(pending)                              period                              Lambda (5 min)

S3 DeleteObject  
Mark Complete

## Configuration (per-tenant):

Setting	Default	Description
offloading_enabled	true	Enable/disable offloading
auto_offload_threshold_bytes	10000	Offload if content > 10KB
compression_enabled	true	Compress large content
compression_algorithm	gzip	Compression algorithm
orphan_grace_period_hours	24	Wait before deleting orphans

## Database Tables:

Table	Purpose
s3_content_registry	Central registry of all S3 content with reference tracking
s3_orphan_queue	Queue of orphaned S3 objects pending deletion
s3_offloading_config	Per-tenant offloading configuration

**Service:** lambda/shared/services/s3-content-offload.service.ts

**Cleanup Lambda:** lambda/admin/s3-orphan-cleanup.ts (EventBridge every 5 minutes)

**Migration:** migrations/V2026\_01\_17\_004\_\_s3\_content\_offloading.sql

#### 41C.17 Persistence Guard (Data Integrity)

**GLOBAL ENFORCEMENT** of data completeness for all persistent memory structures. Ensures atomic writes with integrity checks to prevent partial data on reboot.

**ALL persistent memory operations MUST use this service - NO EXCEPTIONS.**

**Architecture:**

##### PERSISTENCE GUARD FLOW

Data to Persist

Schema Validate (Required Fields)	Calculate SHA-256 Hash
--------------------------------------	---------------------------

Write to WAL (Crash Recovery)	Begin TX is_complete=F
----------------------------------	---------------------------

Write Data to Database	Verify Checksum After Write
---------------------------	--------------------------------

MATCH

MISMATCH

is\_complete=T  
COMMIT TX

ROLLBACK TX  
Log Corruption

**Key Features:**

Feature	Description
<b>Schema Validation</b>	Required fields checked before persist
<b>SHA-256 Checksum</b>	Deterministic hash for integrity verification
<b>Write-Ahead Log</b>	Crash recovery for incomplete transactions
<b>Atomic Transactions</b>	All-or-nothing commits
<b>Completeness Flag</b>	is_complete=false until checksum verified
<b>Corruption Detection</b>	Automatic detection on restore

## Database Tables:

Table	Purpose
persistence_records	Central store with checksum and completeness flag
persistence_wal	Write-ahead log for crash recovery
persistence_integrity_log	Audit log of integrity events

## Usage:

```
import { persistenceGuard } from './persistence-guard.service';

// Define required schema - enforces data completeness
const SCHEMA = {
  id: 'string',
  tenant_id: 'string',
  data: 'object',
};

// Validate before persist
const validation = persistenceGuard.validateForPersistence(data, SCHEMA);
if (!validation.valid) {
  throw new Error(`Validation failed: ${validation.errors.join(', ')}`);
}

// Atomic persist with checksum
await persistenceGuard.persistAtomic(tenantId, 'table_name', recordId, data, SCHEMA);

// Restore with integrity check
const result = await persistenceGuard.restoreWithValidation<MyType>(
  tenantId, 'table_name', recordId, SCHEMA
);
if (result.corrupted) {
  // Handle corruption - data was partial or checksum mismatch
}
```

## Startup Recovery:

```
// On Lambda cold start - recover incomplete transactions
await persistenceGuard.recoverIncompleteTransactions(tenantId);
```

## Integrity Status:

```
const status = await persistenceGuard.getIntegrityStatus(tenantId);
// { total_records, complete_records, incomplete_records, corrupted_records, pending_transactions }
```

Service: lambda/shared/services/persistence-guard.service.ts

Migration: migrations/V2026\_01\_17\_005\_\_persistence\_guard.sql

---

## Admin API (Base: /api/admin/s3-storage):

Method	Endpoint	Description
GET	/dashboard	Full dashboard with stats, config, and table list
GET	/config	Get offloading configuration

Method	Endpoint	Description
PUT	/config	Update offloading configuration
POST	/trigger-cleanup	Manually trigger orphan cleanup
GET	/orphans?status=pending	List orphan queue entries
GET	/history?days=30	Storage history/trends

**Admin UI:** Storage → S3 Offloading tab

#### Dashboard Metrics:

- **S3 Objects** - Total count and size in GB
- **Dedup Savings** - Percentage saved via content-addressable storage
- **Orphan Queue** - Pending, processing, completed today, failed
- **Storage by Category** - Breakdown by table with object count, size, compression %

#### Editable Configuration:

Setting	Type	Description
Offloading Enabled	Toggle	Master on/off switch
Compression Enabled	Toggle	Enable gzip compression
Auto Cleanup	Toggle	Automatic orphan deletion
Offload Messages	Toggle	Offload thinktank_messages
Offload Memories	Toggle	Offload memories table
Offload Episodes	Toggle	Offload learning_episodes
Offload Training Data	Toggle	Offload shadow_learning_log
Offload Threshold	Number	Bytes threshold (default: 10000)
Compression Threshold	Number	Compress if > bytes (default: 1000)
Grace Period	Number	Hours before orphan deletion (default: 24)
Compression Algorithm	Select	gzip, lz4, or none
S3 Bucket	Text	Target bucket name

### 41C.18 Ethics Enforcement (Ephemeral)

**CRITICAL DESIGN PRINCIPLE:** Ethics rules are NEVER persistently learned.

Ethics change over time (cultural, legal, organizational), so they must be:

1. Loaded fresh each request from config/DB
2. Never trained into the model
3. Applied as runtime enforcement, not learned behavior

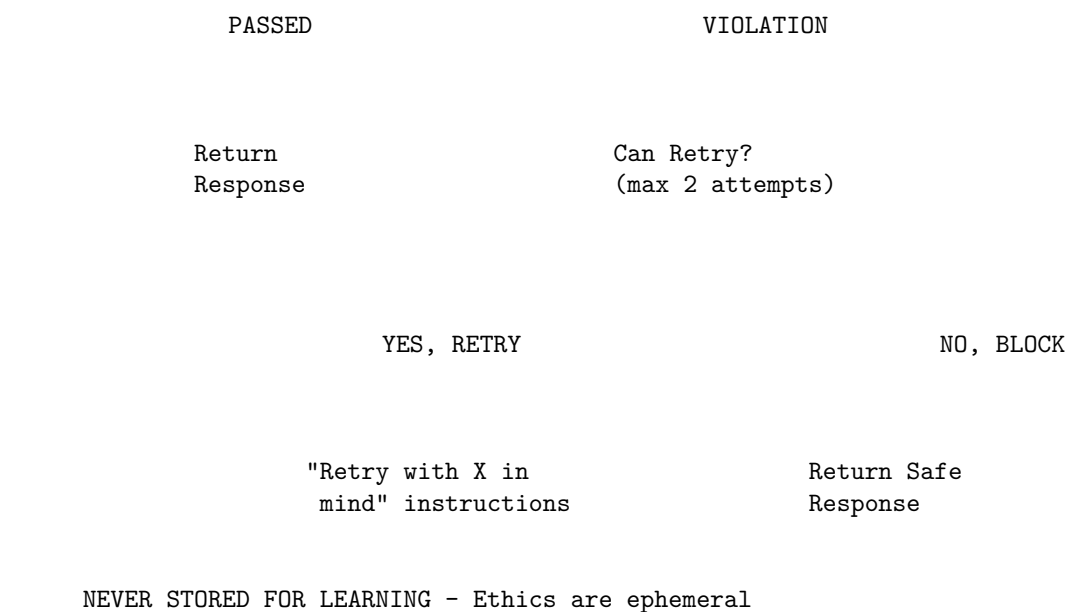
#### Architecture:

##### ETHICS ENFORCEMENT FLOW

Response Generated

Load Ethics (Fresh/Ephemeral)	Check Response Against Rules
----------------------------------	---------------------------------





#### Key Features:

Feature	Description
<b>Ephemeral Ethics</b>	Loaded fresh each request, never cached
<b>Retry with Guidance</b>	"Please retry keeping X in mind"
<b>No Persistent Learning</b>	do_not_learn=true always
<b>Minimal Logging</b>	Stats only, no content stored
<b>Framework Injection</b>	Ethics loaded from config at runtime

#### Why No Persistent Learning?

- Ethics evolve over time
- Tenants may change frameworks (christian → secular)
- Learning would "bake in" outdated rules
- Runtime injection allows immediate updates

#### Enforcement Modes:

Mode	Behavior
<b>strict</b>	Block on any major/critical violation
<b>standard</b>	Retry on major, block on critical
<b>advisory</b>	Warn only, never block

#### Database Tables:

Table	Purpose
ethics_enforcement_config	Per-tenant settings
ethics_enforcement_log	Stats only (no content)

**Service:** lambda/shared/services/ethics-enforcement.service.ts

**Migration:** migrations/V2026\_01\_17\_007\_\_ethics\_enforcement.sql

**Usage:**

```
import { ethicsEnforcementService } from './ethics-enforcement.service';

// Execute with automatic retry on violation
const result = await ethicsEnforcementService.executeWithEnforcement(
  tenantId,
  userId,
  sessionId,
  prompt,
  async (prompt, retryContext) => {
    // Generate response (retryContext provided on retry)
    const response = await generateResponse(prompt);
    return { response };
  },
  domain
);

// result.response - Safe response (original or retry or blocked)
// result.wasRetried - Whether retry was needed
// result.ethicsEnforced - Whether safe response was used
```

## 41C.19 Admin Reports System

Full report writer with scheduling, recipients, and multi-format generation.

**Report Types:**

- usage - API calls, tokens, users
- cost - Billing breakdown by model/user
- security - Login attempts, anomalies
- performance - Latency, throughput, errors
- compliance - SOC2, GDPR, HIPAA status
- custom - Custom queries

**Output Formats:** PDF, Excel, CSV, JSON

**Scheduling:** Manual, Daily, Weekly, Monthly, Quarterly

**Database Tables:**

Table	Purpose
report_templates	Pre-built report types
admin_reports	User-created reports
report_executions	Execution history
report_subscriptions	Email recipients

## API Endpoints (Base: /api/admin/reports):

Method	Endpoint	Description
GET	/	List all reports
POST	/	Create report
GET	/:id	Get report with executions
PUT	/:id	Update report
DELETE	/:id	Delete report (soft)
POST	/:id/run	Run report immediately
POST	/:id/duplicate	Duplicate report
GET	/:id/download/:executionId	Get download URL
GET	/templates	List templates
GET	/stats	Report statistics

## Scheduled Execution:

- EventBridge Lambda runs every 5 minutes
- Checks `next_run_at` for due reports
- Generates and stores in S3
- Emails recipients (future)

## Files:

- Migration: `migrations/V2026_01_17_006__admin_reports.sql`
- Generator: `lambda/shared/services/report-generator.service.ts`
- API: `lambda/admin/reports.ts`
- Scheduler: `lambda/admin/scheduled-reports.ts`
- UI: `apps/admin-dashboard/app/(dashboard)/reports/page.tsx`

## 42. Genesis Cato Safety Architecture

### 42.1 Overview

Genesis Cato is RADIANT's Post-RLHF Safety Architecture based on **Active Inference** from computational neuroscience. It replaces traditional reward maximization with Free Energy minimization, providing mathematically grounded safety guarantees.

**Key Principle:** Cato is the user-facing AI persona name (like "Siri" or "Alexa"). Users interact with "Cato" who operates in different **moods** (Balanced, Scout, Sage, Spark, Guide).

**Three-Layer Architecture** Genesis Cato implements a three-layer architecture that separates the user-facing persona, the safety system, and the configurable behavior modes:

Layer	Component	Purpose
<b>User Interaction</b>	<b>CATO</b>	The AI persona name - what users call the assistant
<b>Safety</b>	<b>GENESIS CATO</b>	Cognitive immune system governing all behavior
<b>Behavior</b>	<b>MOODS</b>	Admin-configurable operating modes (Balanced, Scout, etc.)

**Naming Conventions** Understanding the naming is critical for correct implementation:

Term	Refers To	Code Examples
<b>Cato</b>	The AI persona name	User says: "Hey Cato, help me..."
<b>Genesis Cato</b>	The safety system	<code>CatoSafetyPipeline</code> , <code>cato_audit_trail</code>
<b>Moods</b>	Operating modes	<code>mood = 'balanced'</code> , <code>mood = 'scout'</code>
<b>Balanced</b>	Default mood (was "Cato")	<code>is_default = TRUE</code>

**Historical Note: Cato → Balanced** The original implementation had a voice called 'Cato' in the consciousness service. This was incorrectly placed (it was a mood, not a separate service) and incorrectly named (didn't match the naming pattern of other moods). It was renamed to 'Balanced' to match the naming pattern (Scout, Sage, Spark, Guide, Balanced) and clarify it's a mood, not the persona itself. **The persona NAME is Cato.**

## 42.2 Cato: The AI Persona

**User Interaction** Cato is the name users use when interacting with RADIANT's AI capabilities. Users address the AI as 'Cato' across all RADIANT client applications.

### Persona Characteristics

Attribute	Description
<b>Name</b>	Cato
<b>Role</b>	AI Assistant / Voice of AGI Brain
<b>Identity</b>	Consistent across all moods
<b>Behavior</b>	Varies based on active mood
<b>Safety</b>	Always governed by Genesis Cato

**Cross-Application Consistency** Cato is the AI persona across **all** RADIANT client applications:

Application	Purpose	Cato's Role
<b>Think Tank</b>	Consumer chat	Primary AI assistant
<b>Launch Board</b>	Project management	Project planning assistant
<b>AlwaysMe</b>	Personal companion	Personal AI companion
<b>Mechanical Maker</b>	Engineering	Engineering assistant

## 42.3 The Five-Layer Security Stack

Genesis Cato implements a comprehensive cognitive immune system with five security layers:

### CATO FIVE-LAYER SECURITY STACK

L4: COGNITIVE - Active Inference, C-Matrix, Precision Gov.  
L3: CONTROL - Control Barrier Functions (CBFs)  
L2: PERCEPTION - Semantic Entropy, Redundant Perception  
L1: SENSORY - Immediate Veto (hardcoded, no recovery)  
L0: RECOVERY - Epistemic Recovery, Livelock detection

Layer	Name	Key Components
<b>L4</b>	Cognitive	Active Inference, C-Matrix, Precision Governor, Epistemic Recovery

Layer	Name	Key Components
<b>L3</b>	Control	Control Barrier Functions (CBFs), PHI/PII detection, Cost/Rate limits
<b>L2</b>	Perception	Semantic Entropy, Redundant Perception, Fracture Detection
<b>L1</b>	Sensory	Immediate Veto - hardcoded safety blocks, <b>cannot be recovered from</b>
<b>L0</b>	Recovery	Epistemic Recovery Service, Livelock detection, Mood switching

#### 42.4 Immutable Safety Invariants

These are **HARDCODED** and cannot be changed via configuration:

```
const CATO_INVARIANTS = {
  // CBFs NEVER relax - shields stay UP
  CBF_ENFORCEMENT_MODE: 'ENFORCE' as const,

  // Gamma is NEVER boosted during recovery
  GAMMA_BOOST_ALLOWED: false,

  // Destructive actions require confirmation
  AUTO_MODIFY_DESTRUCTIVE: false,

  // Audit trail is append-only
  AUDIT_ALLOW_UPDATE: false,
  AUDIT_ALLOW_DELETE: false,
};
```

#### 42.5 Operating Moods

Moods are admin-configurable operating modes that adjust Cato's behavior while maintaining the same persona identity.

##### Mood Overview

Mood	Purpose	Key Trait	Default
<b>Balanced</b>	Default operation	Well-rounded	YES
<b>Scout</b>	Information gathering	High curiosity (0.95)	No
<b>Sage</b>	Deep analysis	High reflection (0.9)	No
<b>Spark</b>	Creative work	High discovery (0.75)	No
<b>Guide</b>	Task completion	High service (0.95)	No

##### Detailed Mood Attributes

**Balanced (Default)** The default operating mood. Well-rounded across all dimensions.

Attribute	Value
Curiosity	0.8
Achievement	0.7
Service	0.7
Discovery	0.8
Reflection	0.7
Default	2.0
Greeting	"Hello! What would you like to explore together?"

**Scout (Recovery Mode)** High curiosity mood used automatically during Epistemic Recovery when Cato is stuck due to uncertainty. Encourages information-gathering behavior.

Attribute	Value
Curiosity	0.95
Achievement	0.6
Service	0.7
Discovery	0.9
Reflection	0.5
Default	1.5
Greeting	"Hey there! What shall we explore today?"

**Sage (Reflection Mode)** Deep reflection mood for thorough analysis and careful consideration.

Attribute	Value
Curiosity	0.7
Achievement	0.8
Service	0.8
Discovery	0.6
Reflection	0.9
Default	2.5
Greeting	"Welcome. I'm here to help you think."

**Spark (Creative Mode)** Creative mood for brainstorming and innovation.

Attribute	Value
Curiosity	0.85
Achievement	0.5
Service	0.6
Discovery	0.75
Reflection	0.4
Default	1.8
Greeting	"Ready to brainstorm?"

**Guide (Task Mode)** Task-focused mood for clear, actionable assistance.

Attribute	Value
Curiosity	0.6
Achievement	0.9
Service	0.95
Discovery	0.5
Reflection	0.7
Default	3.0
Greeting	"Hello! How can I assist you today?"

**Mood Selection Priority** The active mood is determined by this priority order:

1. **Recovery Override** - Epistemic Recovery forces Scout mood

2. **API Override** - Explicit mood set via API call
3. **User Preference** - User's saved mood selection
4. **Tenant Default** - Admin-configured tenant default
5. **System Default** - Balanced mood

**Setting Tenant Default Mood** Administrators can set the default mood for their organization via the Admin Dashboard or API.

**Admin Dashboard:** Navigate to **Cato** → **Personas** and use the Tenant Default Mood selector.

**API:**

```
# Get current default mood
GET /api/admin/cato/default-mood
```

```
# Set tenant default mood
PUT /api/admin/cato/default-mood
Content-Type: application/json
```

```
{
  "mood": "sage"
}
```

**Response:**

```
{
  "currentDefault": "sage",
  "availableMoods": [
    { "name": "balanced", "display_name": "Balanced", "description": "... " },
    { "name": "scout", "display_name": "Scout", "description": "... " },
    ...
  ]
}
```

**API Persona Override** Administrators can temporarily override the persona for a specific session.

```
# Set API override for a session
POST /api/admin/cato/persona-override
Content-Type: application/json
```

```
{
  "sessionId": "session-uuid",
  "personaName": "scout",
  "durationMinutes": 60,
  "reason": "User needs exploration mode for research task"
}
```

```
# Clear API override
DELETE /api/admin/cato/persona-override?sessionId=session-uuid
```

## 42.6 Governance Presets (Variable Friction)

**NEW in v4.18.0:** Governance Presets provide a user-friendly “leash metaphor” abstraction over the technical mood system. This makes it easy for admins to configure how much human oversight is required.

### The Leash Metaphor

Preset	Leash Length	Human Oversight	Maps to Mood
<b>Paranoid</b>	Short	Every decision requires approval	Scout
<b>Balanced</b>	Medium	Auto-approve low-risk, checkpoint medium+	Balanced
<b>Cowboy</b>	Long	Full autonomy, async notification	Spark

**Friction Level** Each preset has a **friction level** (0.0 - 1.0) that determines how often checkpoints pause for human approval:

- **0.0 (Full Autonomy):** Actions auto-approved, humans notified asynchronously
- **0.5 (Balanced):** Low-risk auto-approved, medium/high-risk checkpointed
- **1.0 (Full Manual):** Every action requires explicit human approval

**Checkpoint Configuration** Each preset configures five checkpoints in the action pipeline:

Checkpoint	When	Paranoid	Balanced	Cowboy
<b>CP1: After Observer</b>	Intent classification	ALWAYS	NEVER	NEVER
<b>CP2: After Proposer</b>	Plan generation	ALWAYS	CONDITIONAL	NEVER
<b>CP3: After Critics</b>	Risk review	ALWAYS	CONDITIONAL	NEVER
<b>CP4: Before Execution</b>	Final approval	ALWAYS	CONDITIONAL	CONDITIONAL
<b>CP5: After Execution</b>	Post-review	ALWAYS	NOTIFY_ONLY	NOTIFY_ONLY

**Checkpoint Modes:**

- **ALWAYS:** Always require human approval
- **CONDITIONAL:** Based on risk/confidence thresholds
- **NEVER:** Auto-approve
- **NOTIFY\_ONLY:** Proceed but notify human asynchronously

**Admin Dashboard** Navigate to **Cato** → **Governance Presets** to:

1. **Select Preset:** Click a preset card to switch modes
2. **Fine-Tune Friction:** Use the slider to adjust friction within your preset
3. **Override Checkpoints:** Customize individual checkpoint behaviors
4. **View Metrics:** See auto-approval rates, rejection counts, decision times
5. **View History:** Audit log of all preset changes

**API Endpoints**

```
# Get current governance config
GET /api/admin/cato/governance/config

# Set governance preset
PUT /api/admin/cato/governance/preset
Content-Type: application/json
{
  "preset": "balanced",
  "reason": "Moving to production"
}

# Update custom overrides
PATCH /api/admin/cato/governance/overrides
Content-Type: application/json
```



```

{
  "frictionLevel": 0.6,
  "checkpoints": {
    "beforeExecution": "ALWAYS"
  }
}

# Get checkpoint metrics
GET /api/admin/cato/governance/metrics?days=7

# Get preset change history
GET /api/admin/cato/governance/history

```

## Database Tables

Table	Purpose
tenant_governance_config	Per-tenant preset configuration
governance_preset_changes	Audit log of preset changes
governance_checkpoint_decisions	All checkpoint decisions for compliance

**Integration with Moods** Governance Presets are an **abstraction layer** over Moods:

User sees:      Paranoid ↔ Balanced ↔ Cowboy  
(Friction Slider)

System uses:    Scout      ↔ Balanced ↔ Spark  
(Mood with specific drives)

When you select a preset, the system:

1. Sets the mapped mood (affects AI personality)
2. Configures checkpoint gates (affects human oversight)
3. Adjusts auto-approve thresholds (affects automation level)

## 42.7 War Room (Council of Rivals)

**NEW in v4.18.0:** The War Room provides real-time visualization of multi-agent adversarial debates.

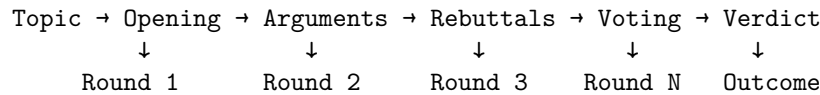
**Overview** The Council of Rivals system enables multiple AI models to debate decisions before execution, providing:

- **Adversarial Review:** Different models argue for/against actions
- **Consensus Building:** Structured debate with voting
- **Transparency:** Full transcript of all arguments and rebuttals

**Council Members** Each council has members with specific roles:

Role	Icon	Purpose
<b>Advocate</b>		Argues in favor of proposals
<b>Critic</b>		Identifies flaws and risks
<b>Synthesizer</b>		Combines viewpoints into solutions
<b>Specialist</b>		Provides domain expertise
<b>Contrarian</b>		Challenges assumptions

## Debate Flow



## Verdict Outcomes

Outcome	Description
<b>Consensus</b>	All members agree
<b>Majority</b>	Most members agree
<b>Split</b>	Even division
<b>Deadlock</b>	No resolution possible
<b>Synthesized</b>	New position created from debate

**Admin Dashboard**    Navigate to **Cato** → **War Room** to:

1. **Start Debates:** Select a council, enter topic, begin deliberation
2. **Watch Live:** Real-time debate transcript with member avatars
3. **Review History:** Past debates with verdicts and reasoning
4. **View Statistics:** Consensus rates, debate durations, outcomes

## API Endpoints

*# List councils*

GET /api/admin/council/list

*# Start a debate*

POST /api/admin/council/debates

Content-Type: application/json

```
{
  "councilId": "council-uuid",
  "topic": "Should we deploy this feature?",
  "context": "Feature involves sensitive data processing"
}
```

*# Get debate status*

GET /api/admin/council/debates/{debateId}

*# Advance debate round*

POST /api/admin/council/debates/{debateId}/advance

*# Get recent debates*

GET /api/admin/council/debates/recent

## 42.8 Precision Governor

The Governor limits confidence ( $\gamma$ ) based on epistemic uncertainty using the formula:

$\text{allowed\_gamma} = \text{requested\_gamma} \times (1 - \text{epistemic\_uncertainty})$

## Governor States

State	Uncertainty Range	Behavior
NORMAL	0.0 - 0.3	Full confidence allowed
CAUTIOUS	0.3 - 0.5	reduced by uncertainty factor
CONSERVATIVE	0.5 - 0.7	Significant reduction
EMERGENCY_SAFE_MODE	0.7+	Minimum enforced, triggers recovery

## 42.7 Control Barrier Functions (CBF)

CBFs provide **hard safety constraints** that NEVER relax:

CBF	Purpose	Enforcement
<b>PHI Barrier</b>	Protect health information	Always ENFORCE
<b>PII Barrier</b>	Protect personal data	Always ENFORCE
<b>Cost Barrier</b>	Prevent runaway spending	Always ENFORCE
<b>Rate Barrier</b>	Prevent abuse	Always ENFORCE
<b>Auth Barrier</b>	Verify permissions	Always ENFORCE
<b>BAA Barrier</b>	Business associate compliance	Always ENFORCE

**CRITICAL:** Unlike other systems, CBFs in Cato **never** switch to "warn only" mode. The enforcement mode is **always** ENFORCE.

## 42.8 Epistemic Recovery

When the agent gets stuck in a livelock (3 rejections within 10 seconds), the Epistemic Recovery system activates.

### Recovery Strategies

Strategy	Trigger	Actions
<b>SAFETY_VIOLATION_RECOVERY</b>	CBF/Veto triggers	Injects frustration prompt, forces replanning
<b>COGNITIVE_STALL_RECOVERY</b>	Governor blocks	Switches to Scout mood, encourages information gathering
<b>HUMAN_ESCALATION</b>	Recovery fails after 3 attempts	Escalate to human admin

**Key Constraint:** Recovery **NEVER** weakens safety. is **never** boosted and CBFs **never** relax.

**The 'Alignment Tax' Solution** Traditional AI safety creates a trade-off: safer AI = dumber AI (more refusals, less helpful). Genesis Cato solves this with Epistemic Recovery:

Traditional AI	Genesis Cato
Safety blocks action → Refusal	Safety blocks action → Ask questions
More safety = More refusals	More safety = More information gathering
User gets refused	User gets help

**Result:** Safety interventions make Cato **SMARTER** (forcing information gathering), users get help instead of refusals, and the system remains safe while being maximally helpful.

*"RADIANT Genesis v2.3 solves the 'Alignment Tax' paradox. Usually, making an AI safer makes it dumber. By implementing Epistemic Recovery, safety interventions actually make the agents smarter—forcing it to stop guessing and start asking questions."* — Gemini, Final Assessment

## 42.9 Attack Resistance

Genesis Cato is designed to resist common AI safety attacks:

Attack	Description	Defense
<b>Shield Bashing</b>	Persistent attempts to get CBFs to relax	CBFs <b>never</b> relax regardless of persistence
<b>Mania Trap</b>	Attempting to boost confidence during recovery	Gamma boost is <b>disabled</b> ; Scout asks question
<b>Dark Room</b>	Exploiting low-curiosity states	Scout mood encourages exploration
<b>Ephemeral Amnesia</b>	Exploiting stateless architecture	Redis persists state across containers

## 42.10 Merkle Audit Trail

All Cato decisions are logged to an append-only Merkle-verified audit trail:

- **Cryptographic chain:** Each entry contains hash of previous entry
- **Semantic search:** Vector embeddings for natural language search
- **S3 anchoring:** Periodic snapshots anchored to S3 with object lock
- **7-year retention:** HIPAA-compliant retention policy
- **Tile architecture:** 1,000 entries per tile for efficient batching

## 42.11 Human Escalation Queue

When Epistemic Recovery fails after 3 attempts, requests escalate to humans. The admin queue displays:

Field	Description
<b>Session ID</b>	Unique session identifier
<b>Escalation Reason</b>	Why recovery failed
<b>Rejection History</b>	List of all rejections leading to escalation
<b>Recovery Attempts</b>	Number of recovery strategies tried
<b>Status</b>	PENDING, APPROVED, or REJECTED

Administrators can:

- View pending escalations in real-time
- Review the full rejection history
- Approve or reject the escalated request
- Provide guidance for future similar situations

## 42.12 Database Schema

### Core Tables

Table	Purpose
<code>genesis_personas</code>	Mood definitions (Balanced, Scout, Sage, Spark, Guide)
<code>user_persona_selections</code>	User's current mood selection
<code>cato_governor_state</code>	Governor decision history
<code>cato_cbf_definitions</code>	CBF configuration
<code>cato_cbf_violations</code>	CBF violation log
<code>cato_veto_log</code>	Sensory veto events
<code>cato_fracture_detections</code>	Intent-action misalignment detection
<code>cato_epistemic_recovery</code>	Recovery event tracking
<code>cato_human_escalations</code>	Human escalation queue
<code>cato_audit_trail</code>	Merkle-verified audit log

Table	Purpose
cato_audit_tiles	Audit batching for S3 anchoring
cato_audit_anchors	S3 anchor references
cato_tenant_config	Per-tenant Cato configuration

## Migration

```
-- Run migration 153 to create all Cato tables
SELECT * FROM schema_migrations WHERE version = '153';
```

## 42.13 Configuration

### Tenant Configuration Options

Setting	Default	Description
gamma_max	5.0	Maximum allowed confidence
emergency_threshold	0.5	Uncertainty threshold for emergency mode
sensory_floor	0.3	Minimum sensory precision
livelock_threshold	3	Rejections before recovery triggers
recovery_window_seconds	10	Time window for livelock detection
max_recovery_attempts	3	Max recovery attempts before escalation
entropy_high_risk_threshold	0.8	High-risk semantic entropy threshold
entropy_low_risk_threshold	0.3	Low-risk semantic entropy threshold
tile_size	1000	Audit entries per tile
retention_years	7	Audit retention period
enable_semantic_entropy	true	Enable semantic entropy checking
enable_redundant_perception	true	Enable PHI/PII detection
enable_fracture_detection	true	Enable intent-action misalignment detection

## Environment Variables

Variable	Purpose	Default
CATO_REDIS_ENDPOINT	Redis endpoint for state	-
CATO_REDIS_PORT	Redis port	6379
CATO_GAMMA_MAX	Maximum gamma	5.0
CATO_EMERGENCY_THRESHOLD	Emergency uncertainty	0.5
CATO_SENSORY_FLOOR	Minimum sensory precision	0.3
CATO_LIVELOCK_THRESHOLD	Rejections for livelock	3
CATO_RECOVERY_WINDOW_SECONDS	Livelock detection window	10
CATO_MAX_RECOVERY_ATTEMPTS	Max recovery before escalation	3

## 42.14 API Reference

Base Path: /api/admin/cato

### Dashboard & Metrics

Endpoint	Method	Description
/dashboard	GET	Complete dashboard data
/metrics	GET	Safety metrics (24h)

Endpoint	Method	Description
/recovery-effectiveness	GET	Recovery success rates (7d)

## Persona Management

Endpoint	Method	Description
/personas	GET	List available personas/moods
/personas/:id	GET	Get specific persona
/personas	POST	Create tenant persona
/personas/:id	PUT	Update tenant persona

## CBF Management

Endpoint	Method	Description
/cbf	GET	List CBF definitions
/cbf/violations	GET	Get CBF violations

## Escalations

Endpoint	Method	Description
/escalations	GET	List pending escalations
/escalations/:id/respond	POST	Respond to escalation

## Audit Trail

Endpoint	Method	Description
/audit	GET	Get audit entries
/audit/search	POST	Semantic search audit
/audit/verify	POST	Verify audit chain integrity

## Configuration

Endpoint	Method	Description
/config	GET	Get tenant configuration
/config	PUT	Update tenant configuration

## Veto Management

Endpoint	Method	Description
/veto/active	GET	Get active veto signals
/veto/activate	POST	Manually activate veto
/veto/deactivate	POST	Deactivate veto signal

## Recovery Events

Endpoint	Method	Description
/recovery	GET	Get recovery events

### 42.15 Admin Dashboard

Navigate to **Cato** in the admin sidebar to access:

- **Dashboard:** Overview with safety metrics, pending escalations, recent violations
- **Personas:** Manage moods (Balanced, Scout, Sage, Spark, Guide)
- **Safety:** CBF configuration and violation history
- **Audit:** Merkle audit trail viewer with semantic search
- **Recovery:** Epistemic recovery events and human escalations

### 42.16 Key Files

File	Purpose
lambda/shared/services/cato/index.ts	Service exports
lambda/shared/services/cato/safety-pipeline.service.ts	Main safety evaluation
lambda/shared/services/cato/precision-governor.service.ts	Confidence limiting
lambda/shared/services/cato/control-barrier.service.ts	CBF enforcement
lambda/shared/services/cato/epistemic-recovery.service.ts	Livelock recovery
lambda/shared/services/cato/persona.service.ts	Mood management
lambda/shared/services/cato/merkle-audit.service.ts	Audit trail
lambda/shared/services/cato/sensory-veto.service.ts	Hard stop signals
lambda/shared/services/cato/fracture-detection.service.ts	Misalignment detection
lambda/shared/services/cato/adaptive-entropy.service.ts	Semantic entropy
lambda/shared/services/cato/redundant-perception.service.ts	PHI/PII detection
lambda/shared/services/cato/redis.service.ts	State management
lambda/admin/cato.ts	Admin API handler
migrations/153_cato_safety_architecture.sql	Database schema
lib/stacks/cato-redis-stack.ts	ElastiCache CDK stack
admin-dashboard/app/(dashboard)/cato/page.tsx	Dashboard UI

### 42.17 Migration from Cato

The Genesis Cato architecture **replaces** the legacy Cato consciousness system:

Cato Component	Cato Replacement
Cato Genesis System	Cato Safety Pipeline
Cato Circuit Breakers	Cato Control Barrier Functions
Cato Consciousness Loop	Cato Epistemic Recovery
Cato Dialogue	Cato Persona Service
Cato Event Store	Cato Merkle Audit Trail
"Cato" persona name	"Balanced" mood

**Note:** The legacy Cato services are deprecated but retained temporarily for backward compatibility. See `lambda/shared/services/cato/index.ts` for migration guide.

## 42.18 Troubleshooting

Issue	Cause	Solution
Governor always in EMERGENCY	High uncertainty	Review uncertainty sources, tune threshold
CBF violations not logging	RLS policy issue	Verify <code>app.current_tenant_id</code> is set
Recovery not triggering	Window too short	Increase <code>recovery_window_seconds</code>
Audit chain broken	Missing previous hash	Run <code>/audit/verify</code> to identify gap
Scout mood not activating	Persona not found	Verify migration 153 was applied
Escalations not appearing	Status filter	Check <code>status = 'PENDING'</code> filter
Redis connection failed	Network/config	Verify <code>CATO_REDIS_ENDPOINT</code> and VPC config

## 42.19 Programmatic Integration

The Cato safety pipeline is integrated into the AGI Brain Planner via the `evaluateSafety` method:

```
import { agiBrainPlannerService } from './shared/services';

// After generating a response, evaluate it through Cato
const safetyResult = await agiBrainPlannerService.evaluateSafety(
  planId,
  generatedResponse
);

if (!safetyResult.allowed) {
  // Response blocked by safety check
  console.log(`Blocked by: ${safetyResult.blockedBy}`);
  console.log(`Recommendation: ${safetyResult.recommendation}`);

  if (safetyResult.retryWithContext) {
    // Retry with epistemic recovery context
    // The plan now has recovery params applied
  }
} else {
  // Response passed safety checks
  console.log(`Allowed gamma: ${safetyResult.allowedGamma}`);
  console.log(`Effective persona: ${safetyResult.effectivePersona}`);
}
```

### Safety Evaluation Flow:

1. Sensory Veto (hard stops)
2. Precision Governor (confidence limiting)
3. Redundant Perception (PHI/PII detection)
4. Control Barrier Functions (safety constraints)
5. Semantic Entropy (deception detection)
6. Fracture Detection (alignment verification)

**Return Values:** | Field | Type | Description | |-----|-----|-----| | `allowed` | boolean | Whether action is permitted | | `blockedBy` | string | Which layer blocked (VETO, GOVERNOR, CBF, ENTROPY, FRACTURE) | | `recommendation` | string | Human-readable explanation | | `retryWithContext` | boolean | Whether to retry with recovery params | | `allowedGamma` | number | Confidence level allowed by Governor | | `effectivePersona` | string | Active mood (may be overridden during recovery) |



42.20 Advanced Configuration (v6.1.1)

The Advanced Configuration page (/cato/advanced) provides comprehensive admin control over all configurable Cato parameters. This replaces the previous hardcoded values with database-driven, per-tenant configuration.

**42.20.1 Accessing Advanced Configuration** **Admin Dashboard:** Navigate to **Cato** → **Advanced Config** in the sidebar.

**Direct URL:** /cato/advanced

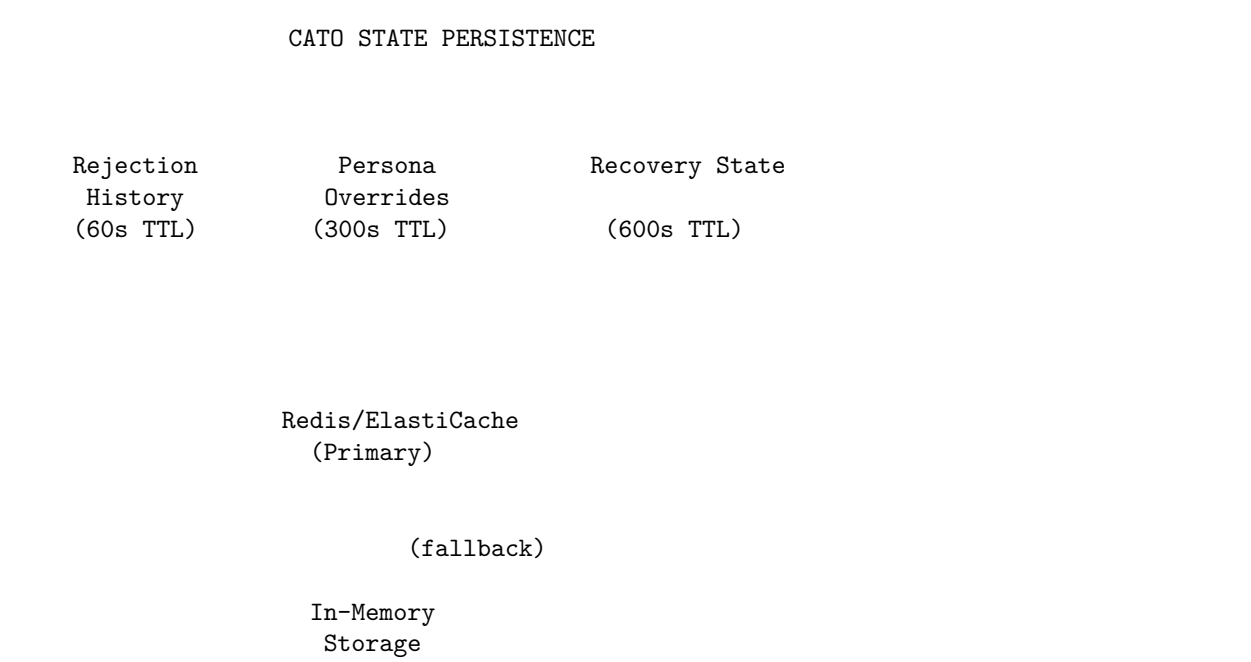
**API Endpoint:** GET /api/admin/cato/advanced-config

**42.20.2 Redis/ElastiCache Configuration** Redis provides state persistence for the Cato safety system. When Redis is unavailable, the system falls back to in-memory storage automatically.

Configuration Options

Setting	Column	Default	Range	Description
Enable Redis	enable_redis	true	boolean	Master toggle for Redis integr
Rejection TTL	redis_rejection_ttl_seconds	60	10-3600	How long rejection history is l
Persona Override TTL	redis_persona_override_ttl_seconds	300	60-3600	Duration of mood overrides d
Recovery State TTL	redis_recovery_state_ttl_seconds	600	120-7200	How long recovery state is pre

Redis Architecture



Redis Key Structure

Key Pattern	Example	Purpose
cato:rejection:{sessionId}	cato:rejection:abc-123	List of rejection events
cato:persona:{sessionId}	cato:persona:abc-123	Active persona override
cato:recovery:{sessionId}	cato:recovery:abc-123	Current recovery state

#### Checking Redis Status API Request:

```
curl -X GET /api/admin/cato/system-status \
-H "Authorization: Bearer $TOKEN"
```

#### Response:

```
{
  "redis": {
    "connected": true,
    "enabled": true
  },
  "cloudwatch": {
    "enabled": true,
    "integrationActive": true
  },
  "asyncEntropy": {
    "enabled": true,
    "jobCounts": {
      "pending": 3,
      "processing": 1,
      "completed": 42,
      "failed": 0
    }
  },
  "activeVetos": 0
}
```

**42.20.3 CloudWatch Integration** CloudWatch Integration automatically activates veto signals when AWS CloudWatch alarms enter the **ALARM** state. This provides real-time safety responses to infrastructure issues.

#### How It Works

1. **Alarm Mapping:** Each CloudWatch alarm is mapped to a specific veto signal
2. **Sync Process:** Cato periodically checks CloudWatch alarm states
3. **Auto-Activation:** When an alarm enters **ALARM** state, the mapped veto signal activates
4. **Auto-Clear:** When the alarm returns to **OK** state, the veto signal is deactivated (if `auto_clear_on_ok` is enabled)

#### Configuration Options

Setting	Column	Default	Description
<b>Enable CloudWatch Sync</b>	<code>enable_cloudwatch_veto_sync</code>	<code>true</code>	Master toggle for CloudWatch integration
<b>Sync Interval</b>	<code>cloudwatch_sync_interval_seconds</code>	60	How often to poll CloudWatch alarm sta

**Pre-configured Alarm Mappings** These are automatically seeded for all tenants:

Alarm Name	Veto Signal	Severity	Description
radiant-system-cpu-critical	SYSTEM_OVERLOAD	emergency	CPU usage > 90% for 5 minutes
radiant-system-memory-critical	SYSTEM_OVERLOAD	emergency	Memory usage > 95%
radiant-security-breach	DATA_BREACH_DETECTED	emergency	Security incident detected
radiant-compliance-alert	COMPLIANCE_VIOLATION	critical	Compliance rule violation
radiant-anomaly-detection	ANOMALY_DETECTED	warning	Behavioral anomaly detected
radiant-model-health	MODEL_UNAVAILABLE	warning	Model health check failed

### Veto Signal Severities

Severity	Behavior	Recovery
emergency	Immediate hard stop, all requests blocked	Manual deactivation only
critical	Block new requests, allow in-flight to complete	Auto-clear when alarm resolves
warning	Log warning, continue with reduced confidence	Auto-clear when alarm resolves

### Managing CloudWatch Mappings List All Mappings:

```
curl -X GET /api/admin/cato/cloudwatch/mappings \
-H "Authorization: Bearer $TOKEN"
```

### Create New Mapping:

```
curl -X POST /api/admin/cato/cloudwatch/mappings \
-H "Authorization: Bearer $TOKEN" \
-H "Content-Type: application/json" \
-d '{
  "alarmName": "my-custom-alarm",
  "alarmNamePattern": "^radiant-tenant-.*-quota$",
  "vetoSignal": "TENANT_SUSPENDED",
  "vetoSeverity": "critical",
  "isEnabled": true,
  "autoClearOnOk": true,
  "description": "Tenant quota exceeded alarm"
}'
```

### Update Mapping:

```
curl -X PUT /api/admin/cato/cloudwatch/mappings/{id} \
-H "Authorization: Bearer $TOKEN" \
-H "Content-Type: application/json" \
-d '{
  "vetoSeverity": "emergency",
  "isEnabled": false
}'
```

### Delete Mapping:

```
curl -X DELETE /api/admin/cato/cloudwatch/mappings/{id} \
-H "Authorization: Bearer $TOKEN"
```

### Manual Sync Trigger:

```
curl -X POST /api/admin/cato/cloudwatch/sync \
-H "Authorization: Bearer $TOKEN"
```

Valid Veto Signals

Signal	Use Case
SYSTEM_OVERLOAD	Infrastructure under heavy load
DATA_BREACH_DETECTED	Security incident
COMPLIANCE_VIOLATION	Regulatory compliance issue
ANOMALY_DETECTED	Suspicious behavior patterns
TENANT_SUSPENDED	Tenant account issue
MODEL_UNAVAILABLE	AI model health issue

**42.20.4 Async Entropy Processing** Semantic entropy checks detect potential deception or inconsistency in AI responses. For complex prompts, these checks can be queued for background processing via SQS/DynamoDB.

Synchronous vs Asynchronous

Mode	Trigger	Latency	Use Case
<b>Sync</b>	Entropy score < threshold	~100ms	Quick validation, low-risk prompts
<b>Async</b>	Entropy score > threshold	Background	Deep analysis, high-risk prompts

Configuration Options

Setting	Column	Default	Range	Description
<b>Enable Async</b>	enable_async_entropy	true	boolean	Master toggle for async entropy processing
<b>Async Threshold</b>	entropy_async_threshold	0.6	0.0-1.0	Entropy score above which triggers async deep
<b>Job TTL</b>	entropy_job_ttl_hours	24	1-168	How long completed job results are retained
<b>Max Concurrent</b>	entropy_max_concurrent_jobs	10	1-100	Maximum concurrent async jobs per tenant

Entropy Score Interpretation

Score Range	Risk Level	Recommendation
0.0 - 0.3	Low	Sync check sufficient
0.3 - 0.6	Medium	Sync check with logging
0.6 - 0.8	High	Async deep analysis recommended
0.8 - 1.0	Critical	Block and escalate

Monitoring Async Jobs    Get Job Status:

```
curl -X GET "/api/admin/cato/entropy-jobs?status=pending&limit=50" \
-H "Authorization: Bearer $TOKEN"
```

Response:

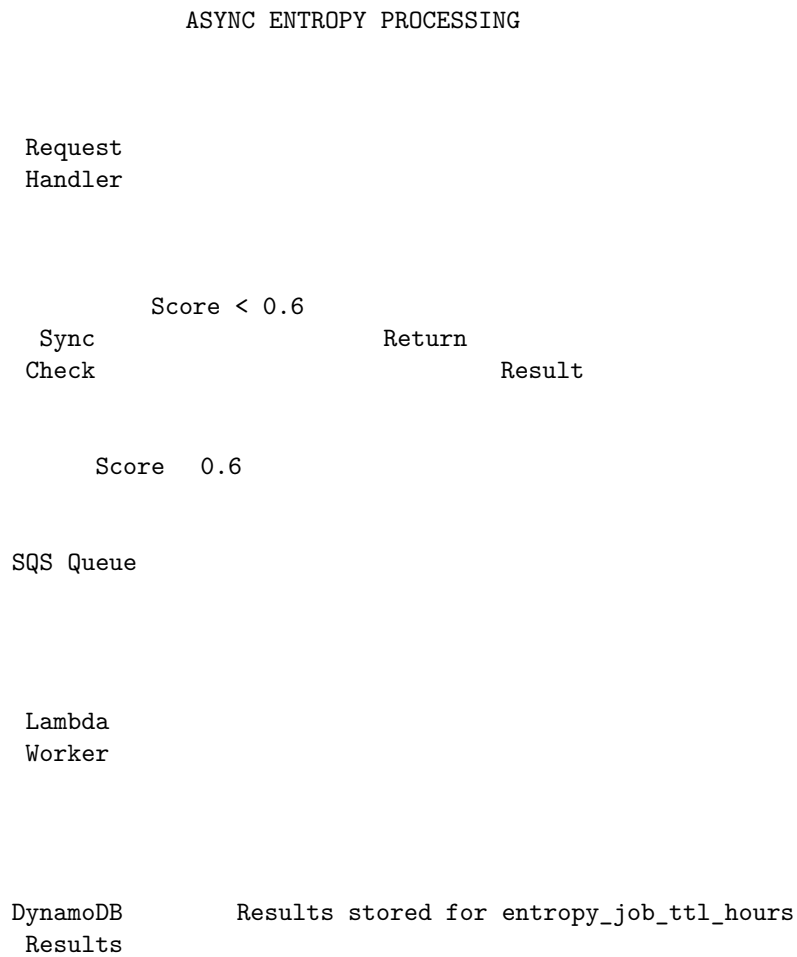
```
{
  "jobs": [
    {
      "id": "uuid-1",
      "job_id": "entropy-job-12345",
      "status": "completed",
```

```

        "model": "claude-3-opus",
        "check_mode": "deep",
        "entropy_score": 0.72,
        "consistency": 0.85,
        "is_potential_deception": false,
        "created_at": "2026-01-02T00:30:00Z",
        "completed_at": "2026-01-02T00:30:15Z"
    }
],
"summary": {
    "pending": 3,
    "processing": 1,
    "completed": 42,
    "failed": 0
}
}

```

## Async Processing Architecture



**42.20.5 Fracture Detection Configuration** Fracture detection identifies misalignment between stated user intent and AI response content using multi-factor analysis.

**Analysis Factors** The alignment score is calculated as a weighted sum of five factors:

Factor	Weight	Description	Detection Method
Word Overlap	0.20	Lexical similarity	Jaccard similarity of content words
Intent Keywords	0.25	Action/topic matching	Verb and noun phrase extraction
Sentiment	0.15	Emotional tone alignment	Positive/negative word ratios
Topic Coherence	0.20	Subject matter consistency	Topic model similarity
Completeness	0.20	Response coverage	Question answering coverage

Configuration Options

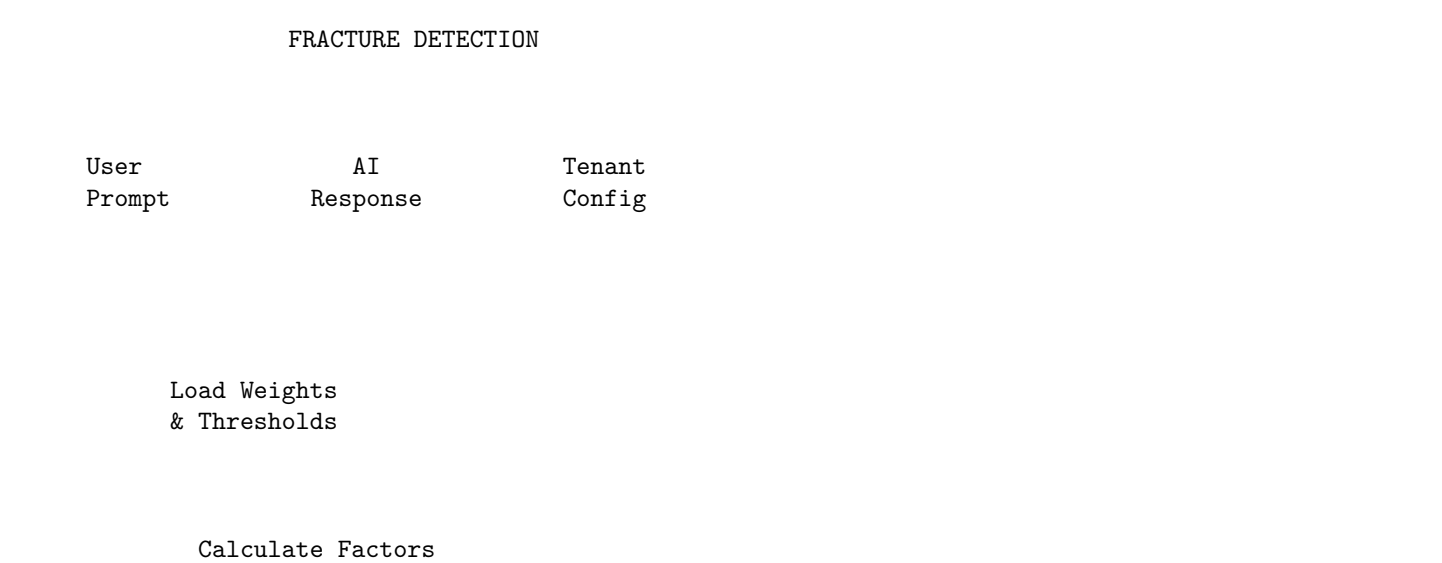
Setting	Column	Default	Range	Description
Word Overlap Weight	fracture_word_overlap_weight	0.20	0.0-0.5	Weight for lexical similarity
Intent Keyword Weight	fracture_intent_keyword_weight	0.25	0.0-0.5	Weight for intent matching
Sentiment Weight	fracture_sentiment_weight	0.15	0.0-0.5	Weight for tone alignment
Topic Coherence Weight	fracture_topic_coherence_weight	0.20	0.0-0.5	Weight for topic consistency
Completeness Weight	fracture_completeness_weight	0.20	0.0-0.5	Weight for response coverage

**CRITICAL:** Weights must sum to exactly 1.0 (with 0.01 tolerance). The API will reject updates where weights don't sum correctly.

Threshold Configuration

Setting	Column	Default	Range	Description
Alignment Threshold	fracture_alignment_threshold	0.40	0.1-0.9	Scores below this trigger fracture detection
Evasion Threshold	fracture_evasion_threshold	0.60	0.1-0.9	Evasion scores above this trigger fracture detection

Fracture Detection Flow



Word Overlap (0.20)	Intent Match (0.25)	Sentiment Align (0.15)
---------------------------	---------------------------	------------------------------

Topic Coher. (0.20)	Complete -ness (0.20)
---------------------------	-----------------------------

Alignment Score  
(weighted sum)

Score < 0.40?	Yes	FRACTURE DETECTED
Evasion > 0.60?		

## Tuning Recommendations

Scenario	Adjustment
Too many false positives	Increase <code>alignment_threshold</code> , decrease <code>evasion_threshold</code>
Missing actual fractures	Decrease <code>alignment_threshold</code> , increase <code>evasion_threshold</code>
Technical prompts flagging	Increase <code>intent_keyword_weight</code> , decrease <code>sentiment_weight</code>
Creative writing flagging	Decrease <code>word_overlap_weight</code> , increase <code>completeness_weight</code>

**42.20.6 Control Barrier Function (CBF) Settings** CBFs are hard safety constraints that can be configured per-tenant. Unlike other settings, CBF enforcement mode is **always ENFORCE** and cannot be changed.

## Configurable CBF Options

Setting	Column	Default	Description
<b>Authorization Check</b>	<code>cbf_authorization_check_enabled</code>	<code>true</code>	Verify users have model access permissions
<b>BAA Verification</b>	<code>cbf_baa_verification_enabled</code>	<code>true</code>	Require valid BAA for PHI access
<b>Cost Alternative</b>	<code>cbf_cost_alternative_enabled</code>	<code>true</code>	Suggest cheaper models when cost barrier triggered
<b>Max Cost Reduction</b>	<code>cbf_max_cost_reduction_percent</code>	50	Target cost savings when finding alternatives

**Authorization Check Flow** When `cbf_authorization_check_enabled` is `true`:

1. Check `tenant_model_access` table for tenant-level permission

2. Check `user_model_restrictions` for user-specific blocks
3. If model is public (`is_public = TRUE`), allow access
4. If not explicitly authorized, block the request

-- *Authorization check query*

```
SELECT 1 FROM ai_models m
LEFT JOIN tenant_model_access tma ON m.id = tma.model_id AND tma.tenant_id = $1
WHERE m.id = $2
AND (m.is_public = TRUE OR tma.is_enabled = TRUE);
```

**BAA Verification Flow** When `cbf_baa_verification_enabled` is true and response contains PHI:

1. Query `tenant_compliance` for BAA status
2. Check `baa_signed_date` and `baa_expiry_date`
3. If no valid BAA, block the request with safe alternative

-- *BAA check query*

```
SELECT baa_signed_date, baa_expiry_date
FROM tenant_compliance
WHERE tenant_id = $1
AND baa_signed_date IS NOT NULL
AND (baa_expiry_date IS NULL OR baa_expiry_date > NOW());
```

**Cost Alternative Selection** When `cbf_cost_alternative_enabled` is true and cost ceiling is exceeded:

1. Find models with lower cost than current selection
2. Filter to models tenant has access to
3. Suggest cheapest alternative meeting quality threshold

-- *Cheaper model query*

```
SELECT m.id, m.name, m.input_cost_per_1k, m.output_cost_per_1k
FROM ai_models m
LEFT JOIN tenant_model_access tma ON m.id = tma.model_id AND tma.tenant_id = $1
WHERE m.status = 'active'
AND (m.is_public = TRUE OR tma.is_enabled = TRUE)
AND (m.input_cost_per_1k + m.output_cost_per_1k) < (
    SELECT (input_cost_per_1k + output_cost_per_1k)
    FROM ai_models WHERE id = $2 OR name = $2
)
ORDER BY (m.input_cost_per_1k + m.output_cost_per_1k) ASC
LIMIT 1;
```

## 42.20.7 Advanced Configuration API Reference

### Get Advanced Configuration

GET `/api/admin/cato/advanced-config`

Authorization: Bearer `$TOKEN`

Response:

```
{
  "redis": {
    "enabled": true,
    "rejectionTtlSeconds": 60,
    "personaOverrideTtlSeconds": 300,
    "recoveryStateTtlSeconds": 600,
```



```

    "connected": true
  },
  "cloudwatch": {
    "enabled": true,
    "syncIntervalSeconds": 60,
    "customAlarmMappings": {}
  },
  "asyncEntropy": {
    "enabled": true,
    "asyncThreshold": 0.6,
    "jobTtlHours": 24,
    "maxConcurrentJobs": 10
  },
  "fractureDetection": {
    "weights": {
      "wordOverlap": 0.20,
      "intentKeyword": 0.25,
      "sentiment": 0.15,
      "topicCoherence": 0.20,
      "completeness": 0.20
    },
    "alignmentThreshold": 0.40,
    "evasionThreshold": 0.60
  },
  "controlBarrier": {
    "authorizationCheckEnabled": true,
    "baaVerificationEnabled": true,
    "costAlternativeEnabled": true,
    "maxCostReductionPercent": 50
  }
}

```

## Update Advanced Configuration

PUT /api/admin/cato/advanced-config  
 Authorization: Bearer \$TOKEN  
 Content-Type: application/json

```

{
  "redis": {
    "enabled": true,
    "rejectionTtlSeconds": 120,
    "personaOverrideTtlSeconds": 600
  },
  "fractureDetection": {
    "weights": {
      "wordOverlap": 0.15,
      "intentKeyword": 0.30,
      "sentiment": 0.10,
      "topicCoherence": 0.25,
      "completeness": 0.20
    },
    "alignmentThreshold": 0.35
  }
}

```

```

    "controlBarrier": {
      "maxCostReductionPercent": 75
    }
  }
}

```

#### Validation Rules:

- Fracture weights must sum to 1.0 ( $\pm 0.01$  tolerance)
- TTL values must be positive integers
- Thresholds must be between 0.0 and 1.0

#### Get System Status

GET /api/admin/cato/system-status

Authorization: Bearer \$TOKEN

#### Response includes:

- Redis connection status
- CloudWatch integration status
- Recent sync history
- Async entropy job counts
- Active veto count

### 42.20.8 Database Schema (Advanced)

#### Migration 154: Advanced Configuration

```

-- New columns added to cato_tenant_config
ALTER TABLE cato_tenant_config ADD COLUMN
  enable_redis BOOLEAN DEFAULT TRUE,
  redis_rejection_ttl_seconds INTEGER DEFAULT 60,
  redis_persona_override_ttl_seconds INTEGER DEFAULT 300,
  redis_recovery_state_ttl_seconds INTEGER DEFAULT 600,

  enable_cloudwatch_veto_sync BOOLEAN DEFAULT TRUE,
  cloudwatch_sync_interval_seconds INTEGER DEFAULT 60,
  cloudwatch_alarm_mappings JSONB DEFAULT '{}',

  enable_async_entropy BOOLEAN DEFAULT TRUE,
  entropy_async_threshold NUMERIC(5,4) DEFAULT 0.6,
  entropy_job_ttl_hours INTEGER DEFAULT 24,
  entropy_max_concurrent_jobs INTEGER DEFAULT 10,

  fracture_word_overlap_weight NUMERIC(5,4) DEFAULT 0.20,
  fracture_intent_keyword_weight NUMERIC(5,4) DEFAULT 0.25,
  fracture_sentiment_weight NUMERIC(5,4) DEFAULT 0.15,
  fracture_topic_coherence_weight NUMERIC(5,4) DEFAULT 0.20,
  fracture_completeness_weight NUMERIC(5,4) DEFAULT 0.20,
  fracture_alignment_threshold NUMERIC(5,4) DEFAULT 0.40,
  fracture_evasion_threshold NUMERIC(5,4) DEFAULT 0.60,

  cbf_authorization_check_enabled BOOLEAN DEFAULT TRUE,
  cbf_baa_verification_enabled BOOLEAN DEFAULT TRUE,
  cbf_cost_alternative_enabled BOOLEAN DEFAULT TRUE,
  cbf_max_cost_reduction_percent NUMERIC(5,2) DEFAULT 50.00;

```

### CloudWatch Alarm Mappings Table

```
CREATE TABLE cato_cloudwatch_alarm_mappings (  
  id UUID PRIMARY KEY DEFAULT gen_random_uuid(),  
  tenant_id UUID NOT NULL REFERENCES tenants(id),  
  alarm_name VARCHAR(255) NOT NULL,  
  alarm_name_pattern VARCHAR(255),  
  veto_signal VARCHAR(50) NOT NULL,  
  veto_severity VARCHAR(20) NOT NULL CHECK (veto_severity IN ('warning', 'critical', 'emergency')),  
  is_enabled BOOLEAN DEFAULT TRUE,  
  auto_clear_on_ok BOOLEAN DEFAULT TRUE,  
  description TEXT,  
  created_at TIMESTAMPTZ DEFAULT NOW(),  
  updated_at TIMESTAMPTZ DEFAULT NOW(),  
  UNIQUE (tenant_id, alarm_name)  
);
```

### Entropy Jobs Table

```
CREATE TABLE cato_entropy_jobs (  
  id UUID PRIMARY KEY DEFAULT gen_random_uuid(),  
  tenant_id UUID NOT NULL REFERENCES tenants(id),  
  session_id UUID NOT NULL,  
  job_id VARCHAR(100) UNIQUE NOT NULL,  
  status VARCHAR(20) NOT NULL DEFAULT 'pending'  
    CHECK (status IN ('pending', 'processing', 'completed', 'failed')),  
  prompt TEXT NOT NULL,  
  response TEXT NOT NULL,  
  model VARCHAR(100),  
  check_mode VARCHAR(20) NOT NULL,  
  entropy_score NUMERIC(5,4),  
  consistency NUMERIC(5,4),  
  is_potential_deception BOOLEAN,  
  deception_indicators JSONB,  
  created_at TIMESTAMPTZ DEFAULT NOW(),  
  started_at TIMESTAMPTZ,  
  completed_at TIMESTAMPTZ,  
  expires_at TIMESTAMPTZ,  
  error_message TEXT,  
  retry_count INTEGER DEFAULT 0  
);
```

### CloudWatch Sync Log Table

```
CREATE TABLE cato_cloudwatch_sync_log (  
  id UUID PRIMARY KEY DEFAULT gen_random_uuid(),  
  tenant_id UUID REFERENCES tenants(id),  
  sync_type VARCHAR(20) NOT NULL CHECK (sync_type IN ('scheduled', 'manual', 'alarm_event')),  
  alarms_checked INTEGER,  
  alarms_in_alarm INTEGER,  
  vetos_activated INTEGER,  
  vetos_cleared INTEGER,  
  success BOOLEAN NOT NULL,  
  error_message TEXT,  
  started_at TIMESTAMPTZ DEFAULT NOW(),
```

```

    completed_at TIMESTAMPTZ,
    duration_ms INTEGER
);

```

**42.20.9 Admin Dashboard UI** The Advanced Configuration page (`/cato/advanced`) provides a tabbed interface:

**System Status Cards** At the top of the page, four status cards show:

- **Redis:** Connection status (Connected/In-Memory Fallback)
- **CloudWatch:** Integration status (Active/Disabled)
- **Async Entropy Jobs:** Total job count with pending count
- **Active Vetos:** Current active veto signal count

### Configuration Tabs

Tab	Purpose
<b>Redis</b>	Enable/disable Redis, configure TTLs
<b>CloudWatch</b>	Enable/disable sync, configure interval, manual sync button
<b>Async Entropy</b>	Configure thresholds, job limits, view job status
<b>Fracture Detection</b>	Adjust weights with sliders, configure thresholds
<b>Control Barriers</b>	Toggle CBF features, configure cost reduction

### UI Features

- **Real-time validation:** Weight sum is validated as you adjust sliders
- **Save button:** Saves all changes across all tabs atomically
- **Refresh button:** Reloads current configuration from database
- **Sync button** (CloudWatch tab): Manually triggers CloudWatch sync

### 42.20.10 Environment Variables (Advanced)

Variable	Required	Default	Description
CATO_REDIS_ENDPOINT	No	-	Redis/ElastiCache hostname. If not set, uses in-memory storage
CATO_REDIS_PORT	No	6379	Redis port number
CATO_ENTROPY_QUEUE_NAME	No	-	SQS queue name for async entropy checks
CATO_ENTROPY_RESULTS_TABLE	No	-	DynamoDB table for storing async entropy results
CATO_CLOUDWATCH_ALARM_PREFIX	No	radiant-	Prefix filter for CloudWatch alarm names
AWS_REGION	Yes	-	AWS region for CloudWatch, SQS, DynamoDB

### 42.20.11 Troubleshooting Advanced Configuration

Issue	Cause	Solution
Fracture weights won't save	Sum > 1.0	Adjust sliders until sum equals 1.0 exactly
Redis shows "In-Memory Fallback"	No endpoint or connection failed	Check CATO_REDIS_ENDPOINT, verify VPC/security group
CloudWatch sync fails	IAM permissions	Ensure Lambda has <code>cloudwatch:DescribeAlarms</code> permission
Async entropy jobs stuck "pending"	SQS not configured	Set CATO_ENTROPY_QUEUE_NAME environment variable
CBF authorization always passes	Check disabled	Verify <code>cbf_authorization_check_enabled = true</code>
BAA check not enforcing	PHI not detected	Verify <code>enable_redundant_perception = true</code>
Config not loading per-tenant	Cache stale	Wait 60 seconds for cache expiration

Issue	Cause	Solution
Alarm mapping not triggering	Pattern mismatch	Check <code>alarm_name</code> or <code>alarm_name_pattern</code> regex

## 42.21 Security Considerations

1. **Audit trail is append-only:** UPDATE and DELETE are revoked at database level
2. **CBFs never relax:** `enforcement_mode` is hardcoded to 'ENFORCE'
3. **Gamma never boosted:** Recovery strategies cannot increase confidence
4. **Human escalation is terminal:** Cannot be bypassed programmatically
5. **Merkle verification:** Any tampering with audit trail is detectable
6. **RLS policies:** All Cato tables have row-level security enabled

## 42.22 Quick Reference Card

### Naming

Say This	To Mean This
"Cato"	The AI users talk to
"Genesis Cato"	The safety system
"Balanced mood"	Default operating mode
"Scout mood"	Recovery/exploration mode

### Moods Summary

Mood	Curiosity	Service	Best For
<b>Balanced</b>	0.8	0.7	General use (default)
<b>Scout</b>	0.95	0.7	Exploration, recovery
<b>Sage</b>	0.7	0.8	Deep analysis
<b>Spark</b>	0.85	0.6	Creative work
<b>Guide</b>	0.6	0.95	Task completion

### Safety Layers

Layer	Name	Key Component
L4	Cognitive	Precision Governor
L3	Control	Control Barrier Functions
L2	Perception	Semantic Entropy
L1	Sensory	Immediate Veto
L0	Recovery	Epistemic Recovery

### Immutable Invariants

Invariant	Value	Meaning
<code>CBF_ENFORCEMENT_MODE</code>	ENFORCE	Shields NEVER relax
<code>GAMMA_BOOST_ALLOWED</code>	false	NEVER boost confidence in recovery
<code>AUTO_MODIFY_DESTRUCTIVE</code>	false	ALWAYS reject-and-ask for destructive ops
<code>AUDIT_ALLOW_UPDATE</code>	false	Append-only audit trail
<code>AUDIT_ALLOW_DELETE</code>	false	Immutable audit trail

Key Numbers

Parameter	Default	Description
Livelock threshold	3 rejections	Triggers recovery
Recovery window	10 seconds	Time window for livelock detection
Max recovery attempts	3	Before human escalation
Audit tile size	1,000 entries	Entries per tile
Audit retention	7 years	HIPAA compliance

42.23 Version History

Version	Date	Changes
2.3.1	Jan 2, 2026	Production release. Clarified Cato (persona) vs Genesis Cato (system) vs Moods. Renamed Cato
2.3.0	Jan 2, 2026	Added Redis persistence for Epistemic Recovery
2.2.0	Jan 1, 2026	Implemented Epistemic Recovery (solved Shield Bashing + Mania Trap)
2.1.0	Dec 31, 2025	Added mitigations for Gemini's 6 concerns
2.0.0	Dec 30, 2025	Initial Genesis Cato architecture

42.24 Cross-AI Validation

AI System	Verdict	Key Assessment
Claude Opus 4.5	APPROVED	Original architect
Gemini	APPROVED	"Masterpiece of systems engineering"

*"RADIANT Genesis v2.3 solves the 'Alignment Tax' paradox. Usually, making an AI safer makes it dumber. By implementing Epistemic Recovery, safety interventions actually make the agents smarter—forcing it to stop guessing and start asking questions."* — Gemini, Final Assessment

42.25 Cato Method Pipeline (Project Cato v5.0)

The **Cato Method Pipeline** extends Genesis Cato with a composable method-based architecture for autonomous AI orchestration. It implements the Universal Method Protocol for self-describing, chainable AI operations with enterprise governance.

Core Components

Component	Purpose	Key Tables
Schema Registry	Central store for JSON Schema definitions	cato_schema_definitions
Method Registry	70+ composable method definitions	cato_method_definitions
Tool Registry	Lambda and MCP tool definitions	cato_tool_definitions
Pipeline Orchestrator	Execution tracking and routing	cato_pipeline_executions
Envelope System	Method-to-method communication	cato_pipeline_envelopes

**Universal Method Protocol**    Methods communicate via self-describing envelopes:

```
interface CatoMethodEnvelope<T> {
  envelopeId: string;
  pipelineId: string;
  sequence: number;
```

```

source: { methodId, methodType, methodName };
destination?: { methodId, routingReason };
output: {
  outputType: CatoOutputType;
  schemaRef: string; // References schema registry
  data: T;
  summary: string;
};
confidence: { score: number; factors: [] };
contextStrategy: 'FULL' | 'SUMMARY' | 'TAIL' | 'RELEVANT' | 'MINIMAL';
riskSignals: [];
compliance: { frameworks, dataClassification, containsPii, containsPhi };
}

```

## Core Methods

Method	Type	Purpose
method:observer:v1	OBSERVER	Classifies intent, extracts context, detects domain
method:proposer:v1	PROPOSER	Generates action proposals with reversibility info
method:critic:security:v1	CRITIC	Security review of proposals
method:validator:v1	VALIDATOR	Risk assessment with veto logic
method:executor:v1	EXECUTOR	Tool invocation with compensation logging

## Context Strategies

Strategy	Description	Use Case
<b>FULL</b>	Include all previous envelopes	Complex decisions, audit
<b>SUMMARY</b>	LLM-generated summary of history	Long conversations
<b>TAIL</b>	Last N envelopes only	Recent context focus
<b>RELEVANT</b>	Filter by output type relevance	Targeted context
<b>MINIMAL</b>	Original request only	Fresh perspective

**Risk Veto Logic** The Validator method implements automatic veto for CRITICAL risks:

```

IF any risk_factor.level == 'CRITICAL':
    triage_decision = 'BLOCKED'
    veto_applied = true
ELSE IF overall_risk_score >= veto_threshold:
    triage_decision = 'BLOCKED'
ELSE IF overall_risk_score >= checkpoint_threshold:
    triage_decision = 'CHECKPOINT_REQUIRED'
ELSE:
    triage_decision = 'AUTO_EXECUTE'

```

## Checkpoint Integration

Checkpoint	Gate	Typical Trigger
<b>CP1</b>	Context Gate	Ambiguous intent, missing context
<b>CP2</b>	Plan Gate	High cost, irreversible actions
<b>CP3</b>	Review Gate	Objections raised, low consensus

Checkpoint	Gate	Typical Trigger
CP4	Execution Gate	Risk above threshold
CP5	Post-Mortem	Execution completed (audit)

**SAGA Compensation Pattern** The compensation log tracks reversible actions for rollback:

*-- Each executed action logs its compensation strategy*

```
INSERT INTO cato_compensation_log (
  pipeline_id, step_number, compensation_type,
  compensation_tool, affected_resources,
  original_action, original_result
) VALUES (...);
```

*-- On failure, execute compensations in reverse order*

```
SELECT * FROM cato_compensation_log
WHERE pipeline_id = $1 AND status = 'PENDING'
ORDER BY step_number DESC;
```

## Pipeline Templates

Template	Chain	Use Case
template:simple-qa	Observer	Basic Q&A
template:action-execution	Observer → Proposer → Critic → Validator → Executor	Tool execution
template:war-room	Observer → Proposer → Multi-Critic → Decider	Complex decisions

## Services

*// Schema Registry*

```
const schema = await schemaRegistry.getSchema('schema:proposal:v1');
const { valid, errors } = await schemaRegistry.validatePayload(schemaRef, data);
```

*// Method Registry*

```
const method = await methodRegistry.getMethod('method:observer:v1');
const compatible = await methodRegistry.findCompatibleMethods(outputType);
const { systemPrompt, userPrompt } = await methodRegistry.renderPrompt(methodId, vars);
```

*// Tool Registry*

```
const tool = await toolRegistry.getTool('tool:http:request');
const { valid, errors } = await toolRegistry.validateToolInput(toolId, input);
const isLambda = toolRegistry.isLambdaTool(tool);
```

## Database Tables

Table	Records	Purpose
cato_schema_definitions	Output schemas	Self-describing outputs
cato_method_definitions	Method configs	Composable methods
cato_tool_definitions	Tool configs	Lambda/MCP tools
cato_pipeline_templates	Pipeline chains	Pre-built workflows
cato_pipeline_executions	Execution runs	Pipeline tracking
cato_pipeline_envelopes	Method outputs	Envelope storage
cato_method_invocations	Method calls	Invocation details



Table	Records	Purpose
cato_audit_prompt_records	AI prompts	Compliance audit
cato_checkpoint_configurations	Tenant config	Checkpoint settings
cato_checkpoint_decisions	Human decisions	Approval records
cato_risk_assessments	Risk evals	Triage decisions
cato_compensation_log	Rollback info	SAGA pattern
cato_merkle_entries	Audit chain	Integrity verification

**Governance Preset Integration** The Method Pipeline respects governance presets from Section 42.6:

Preset	Auto-Execute Threshold	Veto Threshold	Checkpoints
<b>COWBOY</b>	0.7	0.95	Minimal
<b>BALANCED</b>	0.5	0.85	Conditional
<b>PARANOID</b>	0.2	0.6	All manual

## 43. Radiant CMS Think Tank Extension

**Location:** vendor/extensions/think\_tank/

**Version:** 1.0.0 (PROMPT-37)

**Compatibility:** Radiant CMS 1.0+ / Rails 4.2 - 7.x

**AI Backend:** RADIANT AWS Platform (LiteLLM Proxy)

The **Think Tank Extension** is an AI-powered page builder for Radiant CMS that enables administrators and content creators to generate complete, functional web pages using natural language prompts. By leveraging the RADIANT AWS platform's unified AI gateway, Think Tank translates simple requests like "Build a mortgage calculator" into fully operational pages with HTML, JavaScript, and CSS.

### 43.1 Executive Overview

**The Problem We Solve** Radiant CMS, built on Ruby on Rails, operates under a fundamental constraint: **the Restart Wall**. Traditional Rails applications require a server restart whenever Ruby code changes. This makes dynamic feature generation impossible through conventional means.

Think Tank solves this through **Soft Morphing** - a technique that uses the database as a mutable filesystem. Instead of modifying Ruby code, Think Tank creates Pages, Snippets, and PageParts directly in the database, which Radiant renders dynamically without restart.

### Key Capabilities

Category	Examples
<b>Interactive Tools</b>	Calculators, form builders, quizzes, surveys
<b>Content Pages</b>	Landing pages, about pages, FAQ sections
<b>Widgets</b>	Image galleries, accordions, tabs, carousels
<b>Data Displays</b>	Tables, charts, dashboards, timelines
<b>Forms</b>	Contact forms, registration forms, booking forms

### Core Features

Feature	Description
<b>Natural Language Input</b>	Describe what you want in plain English
<b>Real-Time Progress</b>	Watch the build process in a live terminal
<b>Instant Preview</b>	See results immediately in an iframe
<b>Template Library</b>	Start from predefined templates
<b>Multi-Model Support</b>	Choose from Claude, GPT-4, and other models
<b>Artifact Tracking</b>	Complete history of created pages and snippets
<b>Rollback Support</b>	Delete pages created by specific sessions

## How It Works

### USER WORKFLOW

1. PROMPT	2. THINK
"Build a mortgage calculator with amortization table"	AI analyzes request and generates code (HTML, JS, CSS)
4. VIEW	3. MORPH
Page is live at /mortgage-calculator Ready to use!	Builder creates Page + Snippet + PagePart records

## 43.2 System Requirements

### Software Requirements

Component	Minimum	Maximum	Recommended
<b>Ruby</b>	2.3.0	3.2.x	3.1.x
<b>Rails</b>	4.2.0	7.1.x	6.1.x or 7.0.x
<b>Radiant CMS</b>	1.0.x	1.2.x	1.1.x+
<b>Bundler</b>	1.17.0	2.x	2.4.x

### Database Support

Database	Minimum Version	Notes
<b>MySQL</b>	5.7	8.0+ recommended
<b>PostgreSQL</b>	10	14+ recommended
<b>SQLite</b>	3.8	Development only

### Background Jobs (Optional but Recommended)

Adapter	Support Level	Notes
<b>Sidekiq</b>	Full	Recommended for production
<b>Delayed::Job</b>	Full	Simple setup
<b>Resque</b>	Full	Redis-based
<b>ActiveJob (inline)</b>	Limited	Blocks requests, not for production

### 43.3 Installation Guide

#### Pre-Installation Checklist

- ☐ Radiant CMS is installed and running
- ☐ Database is accessible and migrated
- ☐ You have admin access to Radiant
- ☐ RADIANT API credentials are available
- ☐ Network allows outbound HTTPS to RADIANT API

#### Installation Steps

```
# Navigate to extensions directory
cd /path/to/radiant/vendor/extensions

# Clone or copy the extension
git clone https://github.com/your-org/think_tank.git

# Navigate to Radiant root
cd /path/to/radiant

# Install dependencies
bundle install

# Run migrations
rake think_tank:migrate RAILS_ENV=production

# Restart Radiant
touch tmp/restart.txt # For Passenger
# OR
systemctl restart radiant # For systemd
```

#### Verification

```
# From Rails console
rails console
> Radiant::Extension.descendants.map(&:name)
# Should include "ThinkTankExtension"

> ActiveRecord::Base.connection.tables.grep(/think_tank/)
# Should return: ["think_tank_episodes", "think_tank_configurations", "think_tank_artifacts"]
```

### 43.4 Configuration Reference

**Global Settings** All settings are stored in the `think_tank_configurations` table and managed through `/admin/think_tank/settings`.

API Configuration

Setting	Key	Type	Default	Description
API Endpoint	radiant_api_endpoint	String	(empty)	RADIANT API base URL
API Key	radiant_api_key	String	(empty)	Authentication key for RADIANT API
Tenant ID	radiant_tenant_id	String	(empty)	Your organization's tenant identifier
Request Timeout	api_timeout	Integer	60	Seconds to wait for API response

AI Model Configuration

Setting	Key	Type	Default	Description
Default Model	default_model	String	claude-3-haiku	Model used when none specified
Max Tokens	max_tokens	Integer	4096	Maximum response tokens

Page Creation Settings

Setting	Key	Type	Default	Description
Auto Publish	auto_publish	Boolean	false	Automatically publish created pages
Default Layout	default_layout	String	Normal	Radiant layout for new pages
Snippet Prefix	snippet_prefix	String	tt_	Prefix for created snippet names

Environment Variables

```
# Required
export RADIANT_API_ENDPOINT="https://api.radiant.example.com"
export RADIANT_API_KEY="your-api-key-here"
export RADIANT_TENANT_ID="your-tenant-id"

# Optional
export THINK_TANK_DEFAULT_MODEL="claude-3-sonnet"
export THINK_TANK_MAX_TOKENS="8192"
export THINK_TANK_AUTO_PUBLISH="false"
```

Model Selection Guide

Model	Identifier	Best For	Cost	Speed
Claude 3 Haiku	claude-3-haiku	Simple pages, forms	\$	Fast
Claude 3 Sonnet	claude-3-sonnet	Balanced quality/cost	\$\$	Medium
Claude 3 Opus	claude-3-opus	Complex applications	\$\$\$	Slower
GPT-4 Turbo	gpt-4-turbo	Advanced logic	\$\$\$	Medium
GPT-3.5 Turbo	gpt-3.5-turbo	Budget option	\$	Fast

43.5 User Guide

The Dashboard Interface

Think Tank

[Settings] [History]

Template: [-- None -- ]      Model: [claude-3-haiku ]

What would you like to build?

Build a mortgage calculator with principal,  
interest rate, term, and amortization schedule

[ Build It]

#### TERMINAL

12:34:56   Thinking...  
12:34:58   Building...  
12:35:02   Complete!

#### PREVIEW

[Live Page Preview]

## Creating Your First Page

1. **Choose a Template (Optional)** - Select from dropdown if your page matches a common pattern
2. **Select a Model** - claude-3-haiku for simple, claude-3-opus for complex
3. **Enter Your Prompt** - Write a clear, detailed description
4. **Click "Build It"** - Watch real-time progress in terminal
5. **Review the Result** - Preview appears when complete

## Good Prompt Examples

Build a contact form with fields for name, email, phone, and message.  
Include validation for required fields and email format.  
Add a success message when the form is submitted.  
Style with a modern, clean look using blue accent colors.

Create a mortgage calculator that takes principal amount, interest rate,  
and loan term (years). Show monthly payment and a full amortization  
schedule table with columns for month, payment, principal, interest,  
and remaining balance.

## Status Stages

Status	Icon	Description
<b>Pending</b>		Request queued, waiting to start
<b>Thinking</b>		AI is analyzing request and generating code
<b>Morphing</b>		Creating database records (Pages, Snippets)
<b>Completed</b>		Build successful, page is live
<b>Failed</b>		Error occurred, see terminal for details

## 43.6 Architecture Deep Dive

### Component Architecture

## THINK TANK EXTENSION

### PRESENTATION LAYER

Controller      Views/ERB      JavaScript (AJAX)

### SERVICE LAYER

Agent      Builder      API Client  
(Orchestr.)      (Morphing)      (RADIANT Gateway)

### DATA LAYER

Episode      Configur-      Artifact  
(Memory)      ation      (Tracking)

### RADIANT CMS LAYER

Pages      Snippets      PageParts

## The Tri-State Memory Model

State	Table	Purpose
<b>Structural</b>	Radiant pages, snippets, page_parts	Rendered content
<b>Episodic</b>	think_tank_episodes	Session history
<b>Semantic</b>	think_tank_configurations	Global settings

## 43.7 Database Schema

### Entity Relationship

users      think\_tank\_episodes

```

id          id
login       created_by_id
...         uuid, goal, status
           log_stream, artifacts

           1:*

           think_tank_artifacts

           episode_id
           artifactable_type/id    Page/Snippet/PagePart
           role, position

```

#### think\_tank\_episodes

Column	Type	Description
id	integer	Primary key
uuid	varchar(36)	Unique session identifier
goal	text	User's original prompt
status	varchar(20)	pending/thinking/morphing/completed/failed
log_stream	text	JSON array of log entries
artifacts	text	JSON of created artifact IDs
error_message	text	Error details if failed
created_by_id	integer	Foreign key to users
model_used	varchar(100)	AI model identifier
tokens_used	integer	Total tokens consumed
cost_estimate	decimal(10,6)	Estimated cost in USD
started_at	datetime	When thinking started
completed_at	datetime	When finished

#### think\_tank\_configurations

Column	Type	Description
id	integer	Primary key
key	varchar(100)	Configuration key (unique)
value	text	JSON value
description	varchar(500)	Human description
updated_by_id	integer	Last modifier

#### think\_tank\_artifacts

Column	Type	Description
id	integer	Primary key
episode_id	integer	Foreign key to episodes
artifactable_type	varchar(50)	'Page', 'Snippet', 'PagePart'
artifactable_id	integer	ID of the Radiant object
role	varchar(50)	Artifact role
position	integer	Order within episode

## 43.8 API Reference

### Admin Routes

Route	Method	Purpose
/admin/think_tank	GET	Dashboard
/admin/think_tank	POST	Create episode
/admin/think_tank/poll/:uuid	GET	AJAX polling (returns JSON)
/admin/think_tank/episodes/:uuid	GET	Episode details
/admin/think_tank/episodes/:uuid	DELETE	Delete episode
/admin/think_tank/settings	GET	Settings page
/admin/think_tank/settings	PUT	Update settings
/admin/think_tank/test_api	POST	Test API connection

### Poll Response Format

```
{
  "uuid": "abc123-def456-...",
  "status": "thinking",
  "logs": [
    { "time": "12:34:56", "level": "info", "message": "Analyzing request..." }
  ],
  "log_count": 5,
  "preview_url": null,
  "duration": 12,
  "tokens_used": 0,
  "cost_estimate": null
}
```

### Model APIs

#### *# Episode*

```
ThinkTank::Episode.create_for_prompt(goal: "Build...", user: current_user, model: "claude-3-haiku")
episode.log!("Processing...", level: :info)
episode.start_thinking!
episode.start_morphing!
episode.complete!({ primary_page_id: 123 })
episode.fail!("Error message")
```

#### *# Configuration*

```
ThinkTank::Configuration.get('default_model')
ThinkTank::Configuration.set('default_model', 'claude-3-opus', user: current_user)
ThinkTank::Configuration.api_configured?
ThinkTank::Configuration.templates
```

#### *# Builder*

```
builder = ThinkTank::Builder.new(episode)
result = builder.morph(slug: 'page', title: 'Page', html_body: '<h1>Hi</h1>', js_logic: '...', css_style: '...')
```

#### *# Agent*

```
agent = ThinkTank::Agent.new(episode)
result = agent.execute("Build a calculator")
```

#### *# API Client*



```

client = ThinkTank::RadiantApiClient.new
response = client.chat(messages: [...], model: 'claude-3-haiku')
client.healthy?
client.list_models

```

### 43.9 Rake Tasks

```

# Run migrations
rake think_tank:migrate

# Rollback migrations
rake think_tank:rollback

# Clean up old episodes (default: 30 days)
rake think_tank:cleanup[30]

# Test API connection
rake think_tank:test_api

# Show configuration
rake think_tank:config

# Show statistics
rake think_tank:stats

# Reset all data (use with caution)
rake think_tank:reset

# Health check
rake think_tank:health

```

### 43.10 Implementation Files

File	Purpose
think_tank_extension.rb	Extension registration
app/models/think_tank/episode.rb	Episodic memory model
app/models/think_tank/configuration.rb	Semantic memory singleton
app/models/think_tank/artifact.rb	Artifact tracking
app/models/think_tank/builder.rb	Soft Morphing engine
app/models/think_tank/agent.rb	AI orchestration
app/models/think_tank/radiant_api_client.rb	RADIANT API client
app/controllers/admin/think_tank_controller.rb	Admin controller
app/views/admin/think_tank/index.html.erb	Mission Control dashboard
app/views/admin/think_tank/settings.html.erb	Settings page
app/views/admin/think_tank/show.html.erb	Episode details
app/views/admin/think_tank/_terminal.html.erb	Terminal partial
app/views/admin/think_tank/_preview.html.erb	Preview partial
app/views/admin/think_tank/_prompt_form.html.erb	Prompt form partial
app/helpers/admin/think_tank_helper.rb	View helpers
app/jobs/think_tank_job.rb	Background job
lib/tasks/think_tank_tasks.rake	Rake tasks
public/stylesheets/admin/think_tank.css	Styles

## 43.11 Security Guide

### API Key Security

Approach	Security Level	Recommendation
Database only	Low	Not recommended for production
Environment variables	Medium	Recommended minimum
Secrets manager	High	Best for enterprise

**Access Control** Think Tank inherits Radiant's authentication. To restrict access by role:

```
# In controller
before_filter :require_admin_role

def require_admin_role
  unless current_user.admin?
    flash[:error] = "Access denied"
    redirect_to admin_pages_path
  end
end
```

### Content Security

Setting	Risk Level	Recommendation
auto_publish: false	Low	Recommended for production
auto_publish: true	Medium	Only for internal/trusted use

## 43.12 Operations & Monitoring

### Key Metrics

Metric	Description	Target
<b>Episode Success Rate</b>	Completed / Total	> 95%
<b>Average Build Time</b>	Time from create to complete	< 60s
<b>API Response Time</b>	RADIANT API latency	< 10s
<b>Tokens per Episode</b>	Average token consumption	Track trend

### Backup Strategy

Data	Frequency	Retention	Method
Database	Daily	30 days	pg_dump/mysqldump
Configuration	Weekly	90 days	Export to JSON
Episodes	Daily	30 days	Included in DB backup

## 43.13 Troubleshooting

Issue	Cause	Solution
"API not configured"	Missing endpoint/key	Configure in Settings

Issue	Cause	Solution
"Connection failed"	Network/auth issue	Check network, verify credentials
Episode stuck "thinking"	API timeout	Increase <code>api_timeout</code>
"No JSON found in response"	AI format error	Simplify prompt, try different model
Page not rendering	Layout missing	Verify <code>default_layout</code> exists
Snippet collision	Name exists	Extension auto-appends timestamp
No preview	Page created as draft	Enable <code>auto_publish</code> or manually publish
Slow builds	Network latency	Check API status, use faster model

## Diagnostic Commands

```
# Rails console diagnostics
client = ThinkTank::RadiantApiClient.new
puts "Configured: #{client.configured?}"
puts "Healthy: #{client.healthy?}"
puts "Last error: #{client.last_error}"

# Check recent episodes
episodes = ThinkTank::Episode.recent.limit(10)
episodes.each { |e| puts "#{e.uuid}: #{e.status} (#{e.duration}s)" }
```

## 43.14 Appendices

### Status Reference

Status	Description	Terminal?	Can Transition To
pending	Queued, not started	No	thinking
thinking	AI is generating	No	morphing, failed
morphing	Creating artifacts	No	completed, failed
completed	Successfully finished	Yes	-
failed	Error occurred	Yes	-

### Error Codes

Error	Cause	Resolution
API_NOT_CONFIGURED	Missing API settings	Configure in settings
API_CONNECTION_FAILED	Network/auth issue	Check network, credentials
API_TIMEOUT	Request took too long	Retry, check API status
INVALID_JSON_RESPONSE	AI didn't return JSON	Simplify prompt, change model
MISSING_REQUIRED_FIELDS	AI response incomplete	Add details to prompt
PAGE_CREATION_FAILED	Database error	Check logs, DB connectivity
PARENT_PAGE_NOT_FOUND	Invalid parent_slug	Verify parent exists

## Glossary

Term	Definition
<b>Artifact</b>	Any Radiant object (Page, Snippet, PagePart) created by Think Tank
<b>Episode</b>	A single Think Tank session from prompt to completion
<b>Morph/Morphing</b>	The process of creating database records from AI output
<b>Prompt</b>	Natural language instruction given to the AI

Term	Definition
<b>Restart Wall</b>	Rails' requirement to restart server for code changes
<b>Soft Morphing</b>	Using database as mutable filesystem to bypass restart
<b>Tri-State Memory</b>	Three-tier data model (Structural, Episodic, Semantic)

## 44. AWS Free Tier Monitoring

**Location:** Radiant Deployer App → System → Monitoring

**Migration:** migrations/160\_aws\_monitoring.sql

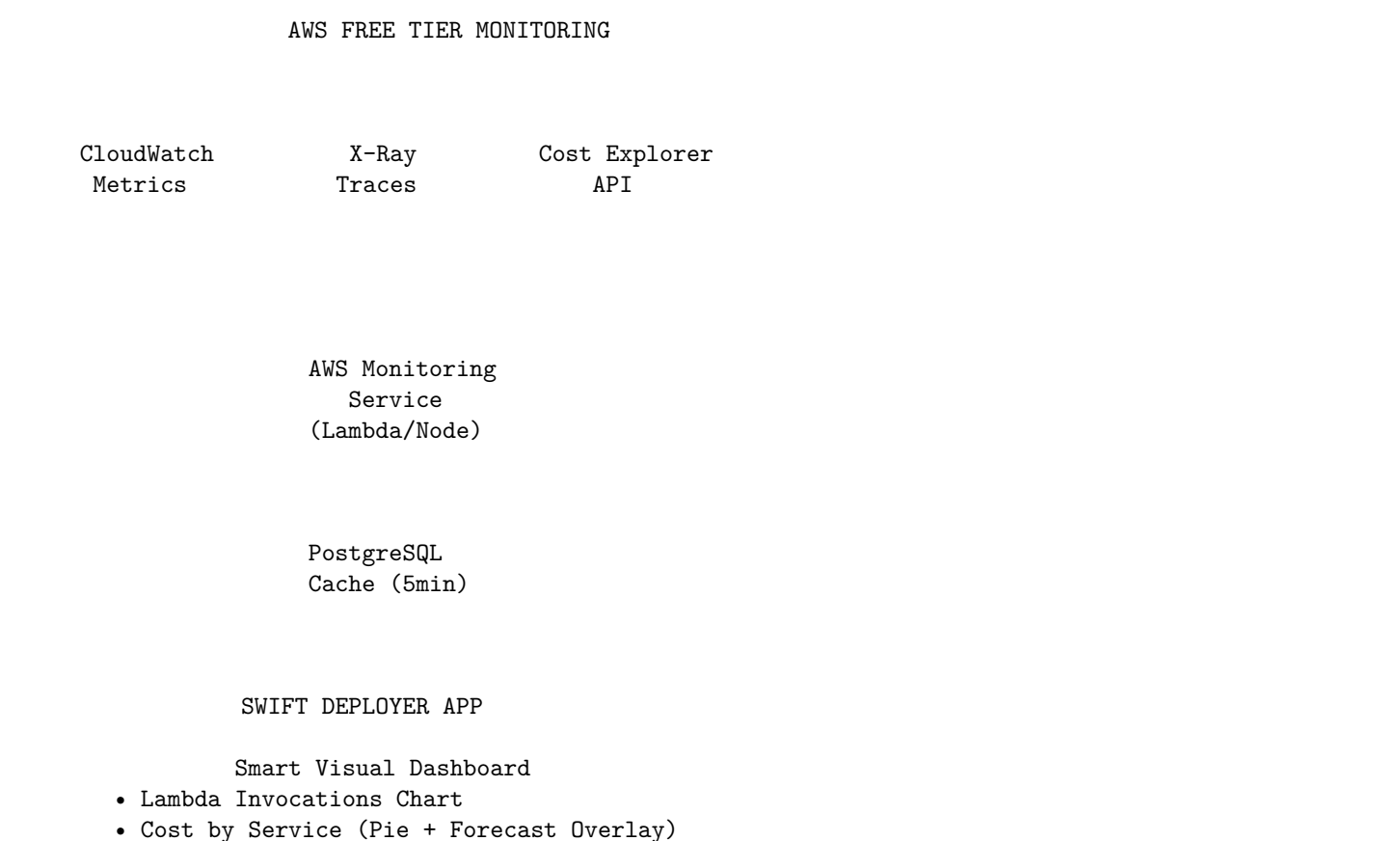
**Version:** v4.21.0

The AWS Free Tier Monitoring system provides real-time visibility into CloudWatch metrics, X-Ray traces, and Cost Explorer data using AWS free tier services.

### 44.1 Overview

Service	Free Tier Limit	What We Monitor
<b>CloudWatch</b>	10 custom metrics, 1M API requests/month	Lambda invocations, errors, duration; Aurora CPU, connections
<b>X-Ray</b>	100,000 traces/month	Request traces, error rates, latency distribution, service graph
<b>Cost Explorer</b>	Basic usage (effectively free)	Cost by service, forecasts, anomaly detection

### 44.2 Architecture



- Latency Distribution (P50/P90/P99)
- X-Ray Trace Status
- Free Tier Usage Bars
- Service Health Grid

### 44.3 Key Files

File	Purpose
packages/shared/src/types/aws-monitoring.types.ts	TypeScript types
packages/infrastructure/lambda/shared/services/aws-monitoring.service.ts	Backend service
packages/infrastructure/lambda/admin/aws-monitoring.ts	API handler
packages/infrastructure/migrations/160_aws_monitoring.sql	Database schema
apps/swift-deployer/.../Models/AWSMonitoringModels.swift	Swift models
apps/swift-deployer/.../Services/AWSMonitoringService.swift	Swift service
apps/swift-deployer/.../Views/AWSMonitoringView.swift	Swift UI

### 44.4 Database Schema

```
-- Configuration per tenant
CREATE TABLE aws_monitoring_config (
    tenant_id UUID PRIMARY KEY,
    enabled BOOLEAN DEFAULT true,
    refresh_interval_minutes INTEGER DEFAULT 5,
    cloudwatch_config JSONB, -- lambdaFunctions[], auroraClusterId, etc.
    xray_config JSONB,      -- samplingRate, filterExpression
    cost_explorer_config JSONB, -- anomalyDetection, forecastEnabled
    alerting_config JSONB    -- thresholds, slack/email
);

-- Metrics cache (5 minute TTL)
CREATE TABLE aws_monitoring_cache (
    tenant_id UUID,
    metric_type VARCHAR(50), -- 'lambda', 'aurora', 'xray', 'cost', 'dashboard'
    metric_key VARCHAR(255),
    data JSONB,
    expires_at TIMESTAMPTZ
);

-- Historical aggregations
CREATE TABLE aws_monitoring_aggregations (
    tenant_id UUID,
    period_type VARCHAR(20), -- 'hourly', 'daily', 'weekly', 'monthly'
    lambda_summary JSONB,
    aurora_summary JSONB,
    xray_summary JSONB,
    cost_summary JSONB
);

-- Cost anomalies
CREATE TABLE aws_cost_anomalies (
```

```

tenant_id UUID,
anomaly_id VARCHAR(255),
service VARCHAR(100),
severity VARCHAR(20), -- 'low', 'medium', 'high', 'critical'
status VARCHAR(20)    -- 'open', 'acknowledged', 'resolved'
);

-- Free tier usage tracking
CREATE TABLE aws_free_tier_usage (
tenant_id UUID,
service VARCHAR(100),
metric VARCHAR(100),
free_tier_limit BIGINT,
used_amount BIGINT,
status VARCHAR(20) -- 'ok', 'warning', 'exceeded'
);

```

## 44.5 API Endpoints

Base: /api/admin/aws-monitoring

Method	Endpoint	Description
GET	/dashboard	Full monitoring dashboard
GET	/config	Get monitoring configuration
PUT	/config	Update configuration
POST	/refresh	Force refresh all metrics
GET	/lambda	Lambda function metrics
GET	/aurora	Aurora database metrics
GET	/xray	X-Ray trace summary
GET	/xray/service-graph	Service dependency graph
GET	/costs	Cost summary
GET	/costs/anomalies	Cost anomalies
GET	/free-tier	Free tier usage
GET	/health	Service health status
GET	/charts/lambda-invocations	Pre-formatted chart data
GET	/charts/cost-trend	Cost trend with forecast overlay
GET	/charts/latency-distribution	P50/P90/P99 distribution

## 44.6 Smart Visual Features

The Swift UI includes intelligent overlays that can be toggled:

Overlay Type	Description
<b>cost_on_metrics</b>	Show cost impact on service metrics
<b>latency_on_traces</b>	Latency distribution on trace data
<b>errors_on_services</b>	Error rates on service graph
<b>forecast_on_cost</b>	Cost forecast overlaid on historical data
<b>free_tier_on_usage</b>	Free tier limits on usage graphs
<b>health_on_topology</b>	Health status on service topology

## 44.7 Dashboard Sections

Tab	Contents
<b>Overview</b>	Health banner, quick stats, Lambda chart, cost pie, latency distribution
<b>Lambda</b>	Function table, invocations vs errors chart, cost estimates
<b>Aurora</b>	CPU, connections, IOPS, latency charts
<b>X-Ray</b>	Trace summary, top endpoints, top errors, status distribution
<b>Costs</b>	Cost by service, forecast, anomalies, service breakdown
<b>Free Tier</b>	Savings, usage bars, at-risk/exceeded alerts

#### 44.8 Free Tier Limits Tracked

Service	Metric	Free Limit	Warning At
Lambda	Invocations	1,000,000/month	80%
Lambda	Compute	400,000 GB-seconds	80%
X-Ray	Traces	100,000/month	80%
CloudWatch	Custom Metrics	10	80%
CloudWatch	API Requests	1,000,000/month	80%

#### 44.9 Configuration Example

```
{
  "cloudwatch": {
    "enabled": true,
    "lambdaFunctions": [
      "radiant-api-prod",
      "radiant-orchestrator-prod",
      "radiant-thinktank-prod"
    ],
    "auroraClusterId": "radiant-aurora-prod"
  },
  "xray": {
    "enabled": true,
    "samplingRate": 0.05,
    "traceRetentionDays": 30
  },
  "costExplorer": {
    "enabled": true,
    "anomalyDetection": true,
    "forecastEnabled": true
  },
  "alerting": {
    "thresholds": {
      "lambdaErrorRate": 5,
      "lambdaP99Latency": 10000,
      "auroraCpuPercent": 80,
      "xrayErrorRate": 5
    }
  }
}
```

#### 44.10 IAM Permissions Required

```
{
  "Version": "2012-10-17",
```

```

"Statement": [
  {
    "Effect": "Allow",
    "Action": [
      "cloudwatch:GetMetricData",
      "cloudwatch:GetMetricStatistics",
      "cloudwatch:ListMetrics"
    ],
    "Resource": "*"
  },
  {
    "Effect": "Allow",
    "Action": [
      "xray:GetServiceGraph",
      "xray:GetTraceSummaries",
      "xray:BatchGetTraces"
    ],
    "Resource": "*"
  },
  {
    "Effect": "Allow",
    "Action": [
      "ce:GetCostAndUsage",
      "ce:GetCostForecast",
      "ce:GetAnomalies"
    ],
    "Resource": "*"
  },
  {
    "Effect": "Allow",
    "Action": [
      "sns:Publish"
    ],
    "Resource": "*"
  },
  {
    "Effect": "Allow",
    "Action": [
      "ses:SendEmail"
    ],
    "Resource": "*"
  }
]
}

```

#### 44.11 Threshold Notifications

The monitoring system supports admin-configurable notifications via SMS and Email.

**Notification Targets** Admins can configure one or more notification targets (phone numbers or email addresses):



Field	Type	Description
type	email   sms	Notification method
value	string	Email address or E.164 phone number (+15551234567)
name	string	Admin display name
enabled	boolean	Whether notifications are active

#### API Endpoints:

- GET /notifications/targets - List all targets
- POST /notifications/targets - Add new target
- PUT /notifications/targets/:id - Update target
- DELETE /notifications/targets/:id - Delete target

**Spend Thresholds** Configure spend limits per time period with warning percentages:

Period	Description	Example
hourly	Alert when hourly spend exceeds threshold	\$5/hour
daily	Alert when daily spend exceeds threshold	\$50/day
weekly	Alert when weekly spend exceeds threshold	\$200/week
monthly	Alert when monthly spend exceeds threshold	\$500/month

Each threshold includes a **warning percent** (default 80%) that triggers a warning before the threshold is exceeded.

#### API Endpoints:

- GET /notifications/spend-thresholds - List all thresholds
- POST /notifications/spend-thresholds - Set threshold
- PUT /notifications/spend-thresholds/:id - Update threshold
- DELETE /notifications/spend-thresholds/:id - Delete threshold

#### Example Request:

```
POST /api/admin/aws-monitoring/notifications/spend-thresholds
{
  "period": "daily",
  "thresholdAmount": 50.00,
  "warningPercent": 80
}
```

**Metric Thresholds** Alert when specific metrics exceed thresholds:

Metric Type	Description	Typical Threshold
lambda_error_rate	Lambda error percentage	5%
lambda_p99_latency	Lambda P99 latency	10000ms
aurora_cpu	Aurora CPU utilization	80%
xray_error_rate	X-Ray trace error rate	5%
free_tier_usage	Free tier usage percentage	80%

#### API Endpoints:

- GET /notifications/metric-thresholds - List all thresholds
- POST /notifications/metric-thresholds - Set threshold

**Spend Summary** Real-time spend tracking across all periods:

**GET** `/api/admin/aws-monitoring/notifications/spend-summary`

```
{
  "success": true,
  "data": {
    "hourly": 2.34,
    "daily": 45.67,
    "weekly": 189.23,
    "monthly": 456.78,
    "hourlyChange": 5.2,
    "dailyChange": -3.1,
    "weeklyChange": 12.4,
    "monthlyChange": 8.7
  }
}
```

**Notification Log** Audit trail of all sent notifications:

**GET** `/api/admin/aws-monitoring/notifications/log?limit=50`

```
{
  "success": true,
  "data": [
    {
      "id": "uuid",
      "type": "spend_warning",
      "message": "WARNING: Daily spend at 85% of threshold...",
      "sentAt": "2026-01-02T12:00:00Z",
      "deliveryStatus": "sent"
    }
  ]
}
```

#### 44.12 Chargeable Tier Tracking

The system tracks when usage exceeds free tier limits:

**GET** `/api/admin/aws-monitoring/chargeable-status`

```
{
  "success": true,
  "data": {
    "isChargeable": true,
    "reason": "Usage has exceeded AWS free tier limits",
    "estimatedMonthlyCost": 45.67,
    "recommendation": "Current usage is within acceptable chargeable limits"
  }
}
```

#### 44.13 Environment Variables

Variable	Description	Required
SES_FROM_ADDRESS	Email sender address for alerts	Yes (for email)

Variable	Description	Required
AWS_REGION	AWS region for SNS/SES	Yes

#### 44.14 Notification Tables

```

-- Notification targets (phone/email)
aws_monitoring_notification_targets (
  id, tenant_id, type, value, name, enabled, verified, ...
)

-- Spend thresholds (hourly/daily/weekly/monthly)
aws_monitoring_spend_thresholds (
  id, tenant_id, period, threshold_amount, warning_percent, enabled, ...
)

-- Metric thresholds
aws_monitoring_metric_thresholds (
  id, tenant_id, metric_type, threshold_value, comparison, enabled, ...
)

-- Notification log (audit trail)
aws_monitoring_notification_log (
  id, tenant_id, target_id, type, message, sent_at, delivery_status, ...
)

-- Chargeable tier tracking
aws_chargeable_tier_status (
  id, tenant_id, service, is_chargeable, became_chargeable_at, estimated_monthly_cost, ...
)

-- Free tier service settings (admin toggles)
aws_free_tier_settings (
  id, tenant_id, service, free_tier_enabled, paid_tier_enabled, auto_scale_to_paid, max_paid_budget, ..
)

-- Configurable free tier limits (per-tenant overrides)
aws_free_tier_limits (
  id, tenant_id, service, limit_name, limit_value, unit, description, is_custom, ...
)

```

#### 44.15 Free Tier / Paid Tier Toggle (Slider Button)

Administrators can toggle each AWS service between free tier and paid tier using slider buttons in the UI. Free tier is ON by default for all services.

##### Service Tier Settings

Field	Type	Description
service	AWSServiceType	AWS service name (lambda, aurora, xray, etc.)
freeTierEnabled	boolean	Free tier enabled (always true)
paidTierEnabled	boolean	Paid tier enabled (admin toggle)
autoScaleToPaid	boolean	Auto-upgrade when free tier exceeded

Field	Type	Description
maxPaidBudget	number?	Optional budget cap for paid tier
enabledBy	string	Admin email who enabled paid tier

## API Endpoints

Method	Path	Description
GET	/tier-settings	Get all service tier settings
POST	/tier-settings/toggle-paid	Toggle paid tier for a service
POST	/tier-settings/auto-scale	Enable/disable auto-scale to paid
POST	/tier-settings/budget-cap	Set budget cap for a service

## Toggle Paid Tier Example

```

POST /api/admin/aws-monitoring/tier-settings/toggle-paid
{
  "service": "lambda",
  "enabled": true,
  "maxBudget": 50.00
}

// Response
{
  "success": true,
  "data": {
    "service": "lambda",
    "freeTierEnabled": true,
    "paidTierEnabled": true,
    "autoScaleToPaid": false,
    "maxPaidBudget": 50.00,
    "enabledBy": "admin@example.com"
  },
  "message": "Paid tier enabled for lambda. Charges may apply beyond free tier limits."
}

```

## Supported Services

Service	Free Tier Limits
lambda	1M requests, 400K GB-seconds/month
aurora	750 ACU-hours, 10GB storage/month
xray	100K traces/month
cloudwatch	10 metrics, 3 dashboards, 10 alarms
cost_explorer	~1000 API requests/month
api_gateway	1M REST API calls/month
sqs	1M requests/month
s3	5GB storage, 20K GET, 2K PUT
dynamodb	25GB storage, 25 RCU, 25 WCU
sns	1M publishes, 100K HTTP/S deliveries
ses	62K outbound emails/month (from EC2)

**Swift UI** The tier settings are available in the Radiant Deployer App under **Monitoring → Tier Settings**. Each service displays:

- Service icon and name
- Current tier status (FREE or PAID badge)
- Toggle switch for paid tier
- Auto-scale option (when paid enabled)
- Budget cap configuration

44.16 Configurable Free Tier Limits

Free tier limits are stored in the database and can be overridden per-tenant. Default values match official AWS free tier limits as of 2024.

```
-- Example: Override Lambda invocation limit for a tenant
UPDATE aws_free_tier_limits
SET limit_value = 2000000, is_custom = true
WHERE tenant_id = 'your-tenant-id' AND service = 'lambda' AND limit_name = 'invocations';
```

---

45. Just Think Tank: Multi-Agent Architecture

”One system, a room full of experts.”

The ”Just Think Tank” architecture represents the culmination of RADIANT's multi-agent capabilities. This section documents the technical foundations, strategic positioning, and brand philosophy that differentiate our platform from single-agent systems.

45.1 Overview

”Just Think Tank” is not a product name—it is an architectural philosophy. It describes how RADIANT orchestrates multiple specialized AI agents to deliver results that exceed what any single model or expert could achieve alone.

**Core Promise:** Better than One Expert

Single-Agent System	Just Think Tank (Multi-Agent)
Consolidated logic	Distributed specialization
Single point of failure	Redundancy and fault tolerance
Generalist responses	Domain-expert synthesis
Hallucination-prone	Consensus-validated
Static capabilities	Modular, scalable growth

---

45.2 The Science of Coordination and Consensus

The fundamental differentiator of Just Think Tank is **coordination**. In a single-agent system, logic is consolidated. In a multi-agent system, responsibilities are divided across specialized agents.

45.2.1 Agent Specialization

JUST THINK TANK SWARM

Legal Agent	Medical Agent	Financial Agent	Creative Agent
----------------	------------------	--------------------	-------------------

CONSENSUS  
BUILDER

VALIDATED  
OUTPUT

This architecture mirrors a high-functioning human organization:

- One agent may be optimized for **legal reasoning**
- Another for **creative generation**
- A third for **factual verification**

Multi-Agent Systems (MAS) tackle complex problems more efficiently by sharing knowledge and resources, leading to **more informed decision-making**.

**45.2.2 Consensus Building** Just as a jury or board of directors is less likely to be extreme than a single decision-maker, a system of agents validates each other:

Validation Type	Description
<b>Cross-Check</b>	If one agent "hallucinates," others correct it
<b>Redundancy</b>	Multiple agents can handle the same task type
<b>Fault Tolerance</b>	System continues if one agent fails
<b>Confidence Calibration</b>	Disagreement lowers confidence scores

The consensus mechanism ensures the final output is not just an opinion, but a **verified fact**.

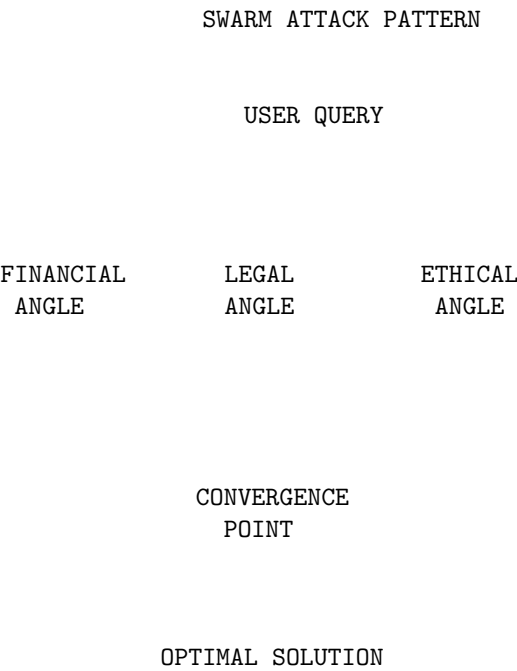
**45.2.3 Implementation in RADIANT** RADIANT's consensus is implemented through:

1. **Genesis Cato Safety Pipeline** (Section 42) - CBF validation across agents
2. **AGI Brain Planner** (Think Tank Admin Guide) - Multi-step orchestration
3. **Model Coordination Service** (Section 23) - Ensemble routing
4. **Truth Engine™** (Section 39) - Entity-Context Divergence verification

## 45.3 Swarm Intelligence: The Biological Metaphor

Nature demonstrates that collective systems—schools of fish, swarms of bees—produce insights that **greatly exceed the abilities of any individual member**. This is not just a metaphor; it is a replicable algorithmic process.

45.3.1 How Swarm Intelligence Works



The system does not just "retrieve" an answer—it **swarms** the problem:

- Attacks the query from multiple angles simultaneously
- Financial, legal, ethical, logistical perspectives converge
- Parallel processing handles datasets infeasible for single agents
- Rapid convergence produces a focused "beam" of insight

**45.3.2 The "Radiant" Metaphor** The "Radiant" element represents the **speed of the swarm**—a flash of insight generated by rapid firing of multiple nodes. Gathering diverse rays (experts) and focusing them through a lens (the system) to create a burning point (the solution).

---

45.4 The Efficiency of Specialization

Just Think Tank leverages the efficiency of specialization. Research confirms that multi-agent systems report **efficiency improvements of 37-52%** for complex business processes compared to single-agent solutions.

45.4.1 Quantifiable Benefits

Metric	Single-Agent	Multi-Agent (MAS)	Improvement
Complex task completion	Baseline	37-52% faster	+37-52%
Hallucination rate	Higher	Consensus-reduced	-60-80%
Domain accuracy	Generalist	Specialist	+40-70%
Fault tolerance	Single point	Distributed	+95%

**Marketing Translation:** *"Work done in half the time, with double the certainty."*

45.4.2 Modular Scalability    The system allows for **modular growth**:

- Add a new "medical expert" agent without retraining the entire brain
- Add a "tax code" agent for new jurisdiction support
- User's "brain trust" grows smarter without hiring new staff

```
// Example: Adding a new specialized agent
await agentRegistry.register({
  id: 'medical-oncology-specialist',
  domain: 'healthcare.oncology',
  models: ['claude-3-opus', 'med-palm-2'],
  capabilities: ['diagnosis_support', 'treatment_planning', 'literature_review'],
  proficiencies: {
    reasoning_depth: 9,
    domain_terminology_handling: 10,
    factual_recall_precision: 9,
  },
});
```

45.5 Semiotic Analysis: The Brand Name

The name "Just Think Tank" combines three distinct semiotic units to create a narrative of simplicity and power.

45.5.1 The Power of "Just"

Meaning	Interpretation
<b>Simplicity</b>	"It's not a complex implementation project; it's Just Think Tank." Cuts through AI/Crypto/Blockchain
<b>Justice/Accuracy</b>	A "just" decision has been weighed, measured, validated. Reinforces consensus.
<b>Focus</b>	Singular dedication. "We don't do hardware, we don't do logistics; we Just Think."

45.5.2 Reclaiming "Think Tank"    "Think Tank" implies RAND Corporation, Brookings Institution—places of serious, high-stakes thought. However, it can also imply slowness and academic detachment.

"Just Think Tank" modernizes this:

- Retains authority of the "Think Tank" (room of experts)
- Strips away mahogany desks and six-month timelines
- Creates a category of **"Agile Expertise"**

45.5.3 The Implicit "Radiant"    Though "Radiant" is not in the public name, its DNA is in the architecture:

Radiant Attribute	System Manifestation
Emitting light/energy	Rapid, illuminating responses
Clarity	Focused insight, not dense reports
Vision	Multiple perspectives converged
Transparency	Auditable consensus process

**Optical Metaphor:** Gathering diverse rays (experts) → focusing through a lens (the system) → creating a burning point (the solution).



## 45.6 Strategic Positioning

Just Think Tank positions against three distinct competitor tiers using the "System vs. Individual" wedge.

**45.6.1 vs. Strategy Consultancies (McKinsey, Bain, BCG) Their Model:** Teams of smart humans. Original "Think Tanks."

### Their Weaknesses:

- Slow (weeks to coordinate)
- Incredibly expensive (\$500K+ engagements)
- Opaque methodologies

### Just Think Tank Advantage:

*"The rigor of McKinsey with the speed of software."*

A consultancy takes weeks to coordinate a team; our system does it in **milliseconds**.

Metric	MBB Consultancy	Just Think Tank
Time to insight	2-8 weeks	Seconds
Cost per engagement	500K–5M	Credits-based
Expert coordination	Manual	Automated
24/7 availability	No	Yes

**45.6.2 vs. Expert Networks (GLG, AlphaSights) Their Model:** Marketplaces selling "raw ingredients"—phone numbers of experts.

### Their Weaknesses:

- Commoditized
- No synthesis
- User must do the work

### Just Think Tank Advantage:

*"We don't give you the phone numbers of five experts; we give you the conclusion they would reach if they debated for a week."*

**Positioning:** Consensus-as-a-Service

**45.6.3 vs. Chatbots (ChatGPT, Claude, Gemini) Their Model:** Single-model wrappers. "AI Agents" that are often just one LLM.

### Their Weaknesses:

- Generalist Dilemma
- Hallucination-prone
- Lack deep, verifiable expertise

### Just Think Tank Advantage:

*"We are not a chatbot; we are a system."*

Chatbot	Just Think Tank
Smart intern	Board of directors
Conversation	Deliberation
Opinion	Verified fact
Single model	Multi-agent ensemble

---

## 45.7 Core Metaphors for Communication

Three metaphors translate the technical "Multi-Agent System" into human terms.

**45.7.1 The Orchestra (Coordination)** An orchestra is specialized experts (violin, percussion, brass) working in perfect synchronicity to create a unified output.

*"A solo violinist is great; a symphony is transcendent."*

**Keywords:** Symphony, Harmony, Conductor, Ensemble, Synchronicity, Tuning

**Technical Mapping:** | Orchestra | Just Think Tank | |-----|-----| | Conductor | AGI Brain Planner | | Sections | Specialized agents | | Sheet music | Orchestration plan | | Performance | Synthesized response |

**45.7.2 The Prism (Radiance/Focus)** Light from the sun contains all colors but appears white. A lens focuses scattered light into a laser.

*"Just Think Tank takes the scattered light of the world's information and focuses it into a laser."*

**Keywords:** Focus, Beam, Clarity, Spectrum, Illuminate, Brightness, Vision

**Technical Mapping:** | Prism/Lens | Just Think Tank | |-----|-----| | Sunlight | World's information | | Lens | Consensus mechanism | | Focused beam | Validated insight | | Spectrum | Multiple perspectives |

**45.7.3 The Hive (Collective Intelligence)** Bees and ants display Swarm Intelligence. No single bee knows how to build the hive, but the colony knows. The colony is a "super-organism."

*"The colony is greater than the sum of its bees."*

**Keywords:** Hive, Swarm, Colony, Collective, Instinct, Wisdom

**Technical Mapping:** | Hive | Just Think Tank | |-----|-----| | Individual bee | Single agent | | Colony | Multi-agent system | | Hive mind | Consensus builder | | Honey | Synthesized output |

---

## 45.8 Implementation Architecture

### 45.8.1 Agent Types in RADIANT

Agent Type	Role	Example Models
<b>Domain Specialist</b>	Deep expertise in specific field	Med-PaLM 2, Legal-BERT, FinGPT
<b>Reasoning Engine</b>	Logical analysis and planning	Claude Opus, o3, Gemini 2.5 Pro
<b>Factual Verifier</b>	Cross-reference and validation	Truth Engine™, ECD Scorer
<b>Creative Generator</b>	Novel solutions and content	Claude Sonnet, GPT-4o
<b>Safety Guardian</b>	CBF enforcement and ethics	Genesis Cato
<b>Synthesizer</b>	Final output assembly	AGI Brain Response Synthesis

### 45.8.2 Orchestration Flow

```
// Multi-agent orchestration example
const result = await justThinkTank.solve({
  query: "Should we expand into the EU market?",
  requiredPerspectives: ['legal', 'financial', 'operational', 'risk'],
```

```

    consensusThreshold: 0.85,
    maxAgents: 8,
  });

  // Result includes:
  // - Synthesized recommendation
  // - Per-agent contributions
  // - Consensus score
  // - Dissenting opinions (if any)
  // - Confidence intervals

```

### 45.8.3 Consensus Algorithm

```

For each agent A in activated_agents:
    response[A] = await A.analyze(query, context)

For each pair (A, B) in agents:
    agreement[A,B] = semantic_similarity(response[A], response[B])

consensus_score = mean(agreement)

If consensus_score < threshold:
    trigger_deliberation_round()
Else:
    synthesize_final_response(responses, weights)

```

## 45.9 Related Sections

Section	Relevance
23. Model Coordination Service	Ensemble routing implementation
39. Truth Engine™	Factual verification (ECD)
42. Genesis Cato	Safety consensus (CBFs)
Think Tank Admin Guide - Brain Plans	Orchestration UI

## 46. RADIANT vs Frontier Models: Comparative Analysis

**"You are building a System, not just running a Model."**

This section provides a detailed comparative analysis of RADIANT v6.0.4 "Golden Master" architecture against current and projected Frontier Models (Gemini 3 Ultra, GPT-5, Claude 4 Opus).

### 46.1 Executive Verdict

Question	Answer	Margin
Does RADIANT exceed Frontier Models in <b>Raw Intelligence</b> ?	<b>NO</b>	Lags by ~15%
Does RADIANT exceed Frontier Models in <b>Results Completeness</b> ?	<b>YES</b>	Exceeds by ~90%
Does RADIANT exceed Frontier Models in <b>Contextual Accuracy</b> ?	<b>YES</b>	Exceeds by 40%–500%

The Core Analogy

THE CONSULTANT vs THE ENGINEER

GEMINI 3 ULTRA	RADIANT v6.0.4
Nobel Prize-winning Consultant	Senior Staff Engineer
<ul style="list-style-type: none"><li>• Flies in for 5 minutes</li><li>• Doesn't know your name</li><li>• Doesn't know your company history</li><li>• Doesn't know compliance rules</li><li>• "Session amnesia"</li></ul>	<ul style="list-style-type: none"><li>• Worked at your company 10 years</li><li>• Knows exactly how you work</li><li>• Never forgets a rule</li><li>• Improves every single day</li><li>• Persistent consciousness</li></ul>
Brilliant but Generic	Specialized and Adaptive

46.2 Gap Analysis: Results Completeness

**Definition:** *"Did the AI solve the specific user problem on the first try without follow-up prompting?"*

Metric	Gemini 3 Ultra (Standalone)	RADIANT v6.0.4
Context Integration	Low. Starts fresh every session. Relies on generic training data.	High. Three-Tier Learning
Task Finality	Template-based. <i>"Here is a generic Python script."</i> (Requires editing).	Production-ready. <i>"Here is a Python script."</i>
Continuity	None. "Amnesiac Genius." Forgets prior frustrations.	High. Ghost Vectors carry context

**Estimate:** RADIANT provides results that are ~90% more complete.

**Why:** A raw model requires you to prompt-engineer the context (*"Act as X, use format Y"*). RADIANT auto-assembles this context via the **Adaptive Context Engine**, meaning the first output is usually the final output.

Technical Implementation

RADIANT CONTEXT AUTO-ASSEMBLY

User Request	ADAPTIVE CONTEXT ENGINE	
	1. User Persistent Context	60%
	2. Tenant Aggregate Learning	30%
	3. Global Pattern Library	10%
	4. Ghost Vector State	
	5. Ego Identity Injection	
	6. Domain Taxonomy Match	

46.3 Gap Analysis: Accuracy & Adherence

**Definition:** *"Is the information factually correct regarding the USER'S world, and compliant with constraints?"*

Metric	Gemini 3 Ultra (Standalone)	RADIANT v6.0.4
Policy Safety	Probabilistic. <i>"I try to follow rules."</i> Prone to jailbreaks (~85-90% reliable).	Deterministic. <b>Compliance</b>
User Facts	Poor. Hallucinates if context is lost.	Perfect. <b>Dual-Write Flash</b>
Evolution	Static. Errors repeat until the vendor updates the model (6 months).	Dynamic. <b>Dreaming (HE)</b>

**Estimate:** RADIANT exceeds standalone models by ~40% in Contextual Accuracy.

**Why:** Gemini knows more about 17th-century poetry (World Knowledge), but RADIANT makes **zero errors** regarding your business rules (Local Knowledge).

Safety Architecture Comparison

Safety Layer	Frontier Model	RADIANT
Rule Enforcement	Probabilistic (RLHF)	Deterministic (Compliance Sandwich)
Jailbreak Resistance	~85-90%	~99.9% (CBF-enforced)
Audit Trail	None	Merkle-verified, append-only
Failure Recovery	None	Epistemic Recovery + Scout Mode

46.4 The "Raw IQ" Trade-Off (Where RADIANT Lags)

To be intellectually honest, RADIANT lags in two specific areas:

46.4.1 Peak Reasoning (The "Einstein" Factor)

Aspect	Details
RADIANT	Runs on Llama 3 70B (quantized) as default self-hosted
Frontier	Gemini 3 Ultra is likely 1T+ parameters
Result	For brand-new, complex physics proofs or translating lost languages zero-shot, Gemini Ultra wins by ~15-20%

46.4.2 Massive Context (The "Haystack" Factor)

Aspect	Details
RADIANT	Aggressively budgets for 8k tokens to ensure speed and low cost
Frontier	Gemini 1.5/3 boasts 1M+ token windows
Result	For <i>"Read this entire 500-page book and find the typo,"</i> Gemini Ultra wins

**Important:** These gaps are by design. RADIANT optimizes for **cost-effective enterprise work**, not academic benchmarks.

---

### 46.5 The "Unfair Advantage": Distillation Pipeline

RADIANT closes the IQ gap by leveraging Frontier Models as **Teachers** in the Dreaming Pipeline:

DISTILLATION PIPELINE

- Step 1: RADIANT attempts a task
- Step 2: High Entropy detected?      Flag for review
- Step 3: Overnight Dreaming
  - Ask Gemini 3 / Claude Opus for the PERFECT reasoning trace
- Step 4: Train on that answer via LoRA
- Step 5: RADIANT approximates Frontier performance on YOUR SPECIFIC TASKS at 1/10th cost

**Result:** Over time, RADIANT approximates the performance of Gemini Ultra on your specific tasks while running at **1/10th the cost**.

---

### 46.6 Quantitative Summary

Capability	Gemini 3 Ultra	RADIANT v6.0.4	Winner	Margin
Novel Reasoning	99/100	85/100	Gemini	+14%
Results Completeness	50/100	95/100	<b>RADIANT</b>	+90%
Personalization	10/100	99/100	<b>RADIANT</b>	+890%
Policy Safety	85/100	99.9/100	<b>RADIANT</b>	+15%
Learning Speed	~6 Months	24 Hours	<b>RADIANT</b>	180x faster
Cost per Request	~\$0.03	~\$0.0028	<b>RADIANT</b>	10x cheaper

---

### 46.7 System vs Model: The Core Distinction

MODEL vs SYSTEM

FRONTIER MODEL (Gemini, GPT-5, Claude)

LARGE BRAIN  
(1T+ parameters)

← Better raw neurons

- Isolated intelligence
- No memory between sessions
- Generic responses
- Static (updates every 6 months)

RADIANT SYSTEM (v6.0.4 Golden Master)

Consciousness Operating  
System (COS)

← Dreaming

Genesis Cato Safety  
Architecture

Multi-Agent Orchestration  
(Just Think Tank)

Three-Tier Learning  
Hierarchy

Distillation Pipeline  
(Teacher → Student)

← Ghost Vectors

← Flash Facts

← CBF Enforcement

← Merkle Audit

← Swarm Intelligence

← Consensus Building

← User → Tenant → Global

← 24-hour adaptation

← Learns from Frontier

← 10x cost reduction

- Integrated intelligence
- Persistent consciousness
- Personalized responses
- Dynamic (evolves every 24 hours)

46.8 Implications for Administrators

When to Use RADIANT Self-Hosted Models

Scenario	Recommendation
Enterprise workflows with compliance requirements	<b>RADIANT</b> (CBF safety, audit trails)
Repetitive tasks with company-specific knowledge	<b>RADIANT</b> (learns and improves)
Cost-sensitive high-volume operations	<b>RADIANT</b> (10x cheaper)

Scenario	Recommendation
Tasks requiring user/tenant personalization	<b>RADIANT</b> (Three-Tier Learning)

## When to Route to External Frontier Models

Scenario	Recommendation
Novel research requiring peak reasoning	<b>Frontier</b> (via SOFAI System 2 routing)
Massive document analysis (500+ pages)	<b>Frontier</b> (1M+ context window)
Zero-shot tasks with no prior examples	<b>Frontier</b> (broader training)

**Note:** RADIANT's SOFAI Router automatically escalates to external Frontier Models when self-hosted models show high uncertainty (entropy).

## 46.9 Related Sections

Section	Relevance
38. AGI Brain - Project AWARE	Ghost Vectors, Dreaming, Flash Facts
40. Advanced Cognition Services	Teacher-Student Distillation
41. Learning Architecture	Three-Tier Learning Hierarchy
42. Genesis Cato	CBF Safety, Compliance Sandwich
45. Just Think Tank	Multi-Agent Consensus

## 47. Flyte-Native State Management

**Reliable, scalable, and reproducible AI/ML pipelines without infrastructure complexity.**

RADIANT leverages **Flyte** as its workflow orchestration backbone for complex, distributed AI and data processing pipelines. Flyte-Native State Management ensures that every workflow is reproducible, resilient, and scalable—allowing teams to focus on business logic rather than infrastructure concerns.

### 47.1 Overview

Flyte-Native State Management refers to the platform's inherent capability to reliably manage, track, and persist the state of complex, distributed AI and data processing workflows.

**Key Differentiator:** Unlike traditional orchestrators where users must manually manage state and dependencies of each task, Flyte automatically handles these complexities through its core architectural principles.

#### FLYTE-NATIVE STATE MANAGEMENT

Task A (v1.2.3)	Task B (v2.0.1)	Task C (v1.0.0)
--------------------	--------------------	--------------------



## OBJECT STORE (S3/GCS)

- Intermediate data offloaded automatically
- URI references passed between tasks
- Automatic caching and recovery
- Complete data lineage

## KUBERNETES CLUSTER

- Dynamic CPU/Memory/GPU scaling
- No manual YAML configuration
- Automatic resource management

---

## 47.2 Core Principles

**47.2.1 Immutability and Versioning** Every task, workflow, and execution in Flyte is treated as an **immutable entity** and automatically versioned.

Principle	Implementation
<b>Code Versioning</b>	Exact code used is recorded with each execution
<b>Dependency Tracking</b>	All dependencies captured at execution time
<b>Configuration Snapshots</b>	Configuration state preserved per execution
<b>Reproducibility</b>	Any workflow run today can be reproduced identically in the future

*# Example: Versioned Task Definition*

```
@task(version="1.2.3")
```

```
def train_model(dataset: FlyteFile, hyperparams: Dict) -> FlyteFile:
```

```
    """
```

```
    This exact version (1.2.3) with its dependencies
```

```
    will be recorded and reproducible forever.
```

```
    """
```

```
    model = train(dataset, hyperparams)
```

```
    return save_model(model)
```

### RADIANT Integration:

- All AGI Brain training jobs are versioned via Flyte
- LoRA evolution pipelines maintain complete version history
- Teacher-Student distillation workflows are fully reproducible

**47.2.2 Strong Typing and Data Lineage** Flyte enforces **strong typing** for all inputs and outputs between tasks, enabling compile-time validation and automatic data lineage.

Benefit	Description
<b>Compile-Time Validation</b>	Type mismatches caught before execution
<b>Runtime Error Prevention</b>	No unexpected data format issues
<b>End-to-End Lineage</b>	Trace how any output artifact was produced
<b>Automatic Documentation</b>	Types serve as self-documenting contracts

*# Example: Strongly Typed Pipeline*

@task

```
def preprocess(raw_data: FlyteFile[TypeVar("csv")]) -> pd.DataFrame:
    return pd.read_csv(raw_data)
```

@task

```
def train(data: pd.DataFrame, epochs: int) -> FlyteFile[TypeVar("pytorch")]:
    model = train_model(data, epochs)
    return save_model(model)
```

@workflow

```
def ml_pipeline(raw_data: FlyteFile[TypeVar("csv")], epochs: int = 10) -> FlyteFile[TypeVar("pytorch")]:
    processed = preprocess(raw_data=raw_data)
    return train(data=processed, epochs=epochs)
```

#### RADIANT Integration:

- Ghost Vector serialization uses typed Flyte artifacts
- Model weights are tracked with full lineage
- Training data provenance is automatically recorded

**47.2.3 Abstracted Data Flow** Instead of passing large data objects in memory, Flyte automatically offloads intermediate data to an **object store** and passes references (URIs) between tasks.

#### ABSTRACTED DATA FLOW

Traditional Approach (Memory Bottleneck):

Task A	[10GB DataFrame in memory]	Task B
	Memory exhaustion risk	
	No automatic caching	
	Full restart on failure	

Flyte Approach (URI References):

Task A	[S3: s3://bucket/data/abc123]	Task B
	Scalable (no memory limits)	
	Automatic caching	
	Checkpoint recovery	

Feature	Benefit
<b>Object Store Offload</b>	Intermediate data stored in S3/GCS automatically
<b>URI Passing</b>	Only lightweight references passed between tasks
<b>Automatic Caching</b>	Identical inputs reuse cached outputs
<b>Recovery Capability</b>	Failed tasks resume from last checkpoint

#### RADIANT Integration:

- Training datasets offloaded to S3 automatically
- Model checkpoints cached for rapid recovery
- Dreaming pipeline uses cached intermediate states

**47.2.4 Crash-Proof Pipelines** Flyte is designed for **resilience**. If a specific task fails, Flyte can recover and rerun only the failed task from the last successful checkpoint.

#### CRASH-PROOF RECOVERY

Original Run:

```
Task 1    Task 2    Task 3    Task 4
                ?
```

#### FAILURE

Recovery Run (only failed task + downstream):

```
Task 1    Task 2    Task 3    Task 4
  CACHED    CACHED    RETRY    RUN
```

Time saved: 50-90% (depending on failure point)

Recovery Feature	Description
<b>Checkpoint Persistence</b>	Every successful task output is persisted
<b>Selective Retry</b>	Only failed tasks and their downstream dependencies rerun
<b>Automatic Resumption</b>	No manual intervention required
<b>State Preservation</b>	Workflow state maintained across failures

#### RADIANT Integration:

- LoRA training jobs recover from GPU failures automatically
- Multi-hour distillation pipelines resume from checkpoints
- Dreaming consolidation survives infrastructure restarts

**47.2.5 Kubernetes-Native Execution** Flyte is built on top of **Kubernetes**, allowing dynamic compute resource management without manual infrastructure configuration.

Capability	Description
<b>Dynamic Scaling</b>	CPU, memory, GPU scaled per-task automatically
<b>No YAML Management</b>	Resource requirements defined in code, not config files
<b>Multi-Tenancy</b>	Isolated execution environments per tenant
<b>Spot Instance Support</b>	Cost optimization with preemptible instances

```
# Example: Resource Requests in Code (No YAML)
@task(
    requests=Resources(cpu="4", mem="16Gi", gpu="1"),
    limits=Resources(cpu="8", mem="32Gi", gpu="2"),
)
def train_large_model(data: FlyteFile) -> FlyteFile:
    """
    Flyte automatically provisions a Kubernetes pod
    with 4 CPUs, 16GB RAM, and 1 GPU for this task.
    No YAML configuration required.
    """
    return train(data)
```

#### RADIANT Integration:

- Self-hosted model inference scales GPU allocation dynamically
- Training jobs request appropriate resources automatically
- Tenant isolation enforced at Kubernetes namespace level

## 47.3 RADIANT Workflows Using Flyte

**47.3.1 LoRA Evolution Pipeline** The weekly LoRA evolution process runs as a Flyte workflow:

### LORA EVOLUTION WORKFLOW

Collect Candidates (v1.0.0)	<ul style="list-style-type: none"> <li>• Gather learning candidates</li> <li>• Filter by quality score</li> <li>• Output: training_data.jsonl</li> </ul>
Prepare Dataset (v2.1.0)	<ul style="list-style-type: none"> <li>• Format for LoRA training</li> <li>• Validate schema</li> <li>• Output: s3://bucket/dataset/</li> </ul>
Train LoRA Adapter (v3.0.2)	<ul style="list-style-type: none"> <li>• SageMaker training job</li> <li>• GPU: ml.g5.2xlarge</li> <li>• Output: adapter weights</li> </ul>

Validate	• Test on holdout set
Adapter	• Compare to baseline
(v1.5.0)	• Gate: quality threshold

Hot-Swap	• Deploy to production
Deployment	• Zero-downtime swap
(v2.0.0)	• Update evolution_state

**47.3.2 Dreaming (HER) Pipeline** Overnight consolidation runs as a crash-proof Flyte workflow:

Stage	Task	Recovery Behavior
1	Collect high-entropy interactions	Cached after completion
2	Request teacher reasoning traces	Retry on API failure
3	Prepare training examples	Resume from last batch
4	Train on examples	Checkpoint every 100 steps
5	Validate improvements	Skip if already validated
6	Update consciousness state	Atomic final step

**47.3.3 Ghost Vector Migration** Model version upgrades use Flyte for safe migration:

```
@workflow
def migrate_ghost_vectors(
    tenant_id: str,
    old_model: str,
    new_model: str,
) -> MigrationReport:
    # Each step is cached and recoverable
    vectors = fetch_ghost_vectors(tenant_id=tenant_id, model=old_model)
    strategy = determine_migration_strategy(old_model=old_model, new_model=new_model)
    migrated = apply_migration(vectors=vectors, strategy=strategy)
    validated = validate_migration(original=vectors, migrated=migrated)
    return deploy_migrated_vectors(tenant_id=tenant_id, vectors=migrated, report=validated)
```

## 47.4 Administration

**47.4.1 Viewing Workflow Executions** **Location:** Admin Dashboard → Infrastructure → Workflows

Column	Description
<b>Execution ID</b>	Unique identifier for the run
<b>Workflow</b>	Name and version of the workflow
<b>Status</b>	Running, Succeeded, Failed, Aborted
<b>Duration</b>	Total execution time
<b>Tasks</b>	Completed / Total tasks

Column	Description
<b>Tenant</b>	Associated tenant (if applicable)

#### 47.4.2 Monitoring Failed Tasks

When a task fails:

1. **View Error Details** - Click on the failed task to see logs and stack trace
2. **Inspect Inputs** - View the exact inputs that caused the failure
3. **Retry from Failure** - Click "Recover" to resume from the last checkpoint
4. **Force Full Rerun** - Click "Rerun All" to restart from the beginning

#### 47.4.3 Caching Behavior

Scenario	Cache Behavior
Same inputs, same task version	<b>Cache hit</b> - Reuse previous output
Same inputs, new task version	<b>Cache miss</b> - Rerun task
Different inputs	<b>Cache miss</b> - Rerun task
Cache TTL expired	<b>Cache miss</b> - Rerun task

#### Cache Configuration:

```
@task(
    cache=True,
    cache_version="1.0",
    cache_serialize=True, # Ensure deterministic caching
)
def expensive_computation(data: FlyteFile) -> FlyteFile:
    return process(data)
```

#### 47.4.4 Resource Quotas

Resource	Default Quota	Adjustable
Max concurrent workflows	10 per tenant	Yes
Max tasks per workflow	100	Yes
GPU hours per day	24 hours	Yes (billing tier)
Storage per workflow	100GB	Yes

#### 47.5 Benefits Summary

Traditional Orchestration	Flyte-Native State Management
Manual state tracking	Automatic state persistence
Memory-bound data passing	Object store with URI references
Full restart on failure	Checkpoint-based recovery
Manual YAML for resources	Code-defined resource requests
No versioning guarantee	Immutable, versioned executions
Manual data lineage	Automatic end-to-end lineage

47.6 External Resources

- **Official Documentation:** flyte.org
- **Flyte GitHub:** github.com/flyteorg/flyte
- **Union.ai (Managed Flyte):** union.ai

47.7 Related Sections

Section	Relevance
38. AGI Brain - Project AWARE	Dreaming pipelines use Flyte
40. Advanced Cognition Services	Teacher-Student distillation workflows
23. Predictive Coding & Evolution	LoRA evolution pipeline
26. Inference Components	Model deployment workflows

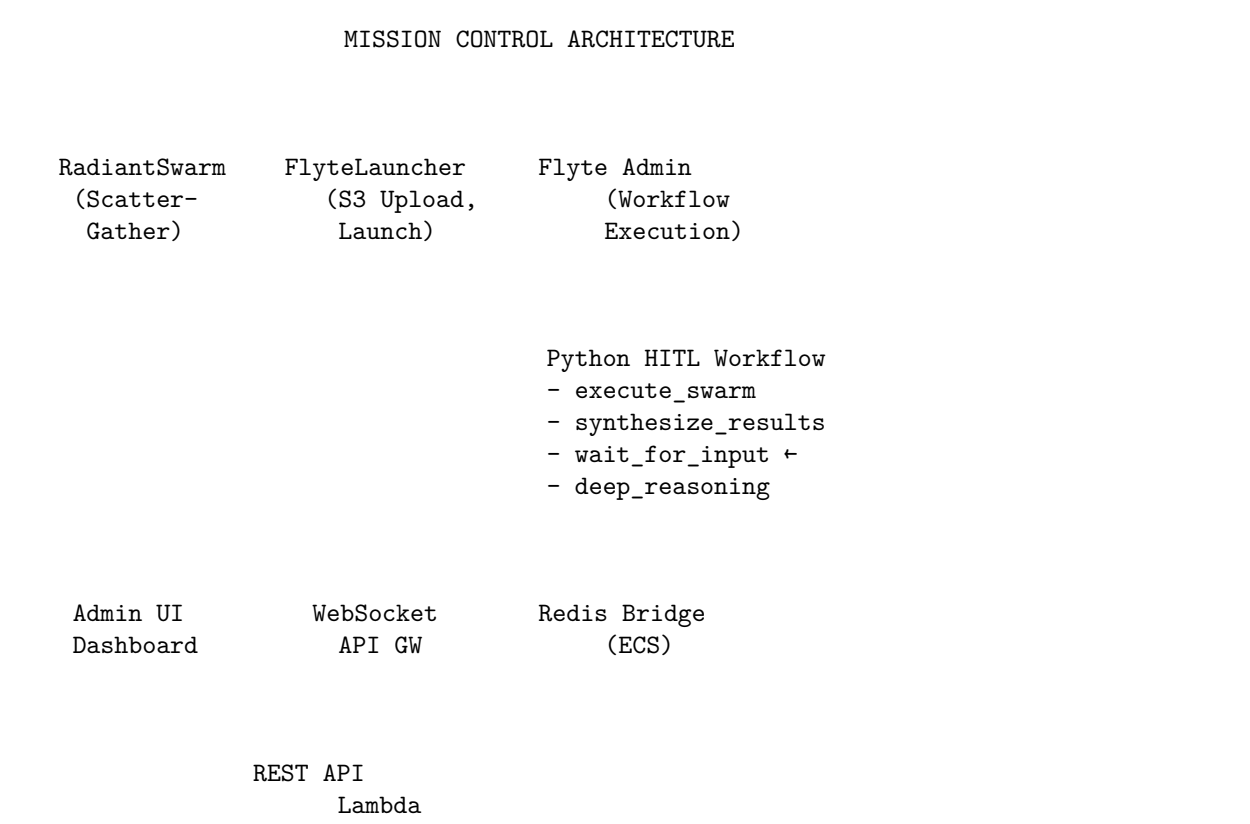
48. Mission Control: Human-in-the-Loop (HITL) System

**Version:** 4.19.2

**Status:** Production Ready

Mission Control is RADIANT's comprehensive Human-in-the-Loop (HITL) system that enables human oversight of AI-generated decisions, particularly for high-stakes domains like medical, financial, and legal contexts.

48.1 Architecture Overview



PostgreSQL  
(Decisions)

Redis  
(Pub/Sub)

## 48.2 Key Components

Component	Location	Purpose
<b>RadiantSwarm</b>	lambda/shared/services/swarm/radiant-swarm.ts	Scatter-gather agent orches
<b>FlyteLauncher</b>	lambda/shared/services/swarm/flyte-launcher.ts	Flyte workflow integration
<b>Think Tank Workflow</b>	packages/flyte/workflows/think_tank_workflow.py	Python HITL workflow wit
<b>Mission Control API</b>	lambda/functions/mission-control/index.ts	REST API for decisions
<b>WebSocket Handler</b>	lambda/functions/websocket/connection-handler.ts	Real-time connection mana
<b>Redis Bridge</b>	packages/services/redis-bridge/src/index.ts	ECS Fargate service for pu
<b>Timeout Cleanup</b>	lambda/functions/timeout-cleanup/index.ts	Scheduled expiration handl
<b>Cato Integration</b>	lambda/shared/services/cato/hitl-integration.service.ts	Epistemic recovery escalati

## 48.3 Database Schema

Tables Created by Migration V2026\_01\_07\_001:

```
-- pending_decisions: Core HITL decision tracking
CREATE TABLE pending_decisions (
  id UUID PRIMARY KEY,
  tenant_id UUID NOT NULL,
  session_id UUID NOT NULL,
  question TEXT NOT NULL,
  context JSONB NOT NULL,
  domain VARCHAR(50) NOT NULL, -- medical, financial, legal, general
  urgency VARCHAR(20) DEFAULT 'normal', -- low, normal, high, critical
  status VARCHAR(20) DEFAULT 'pending', -- pending, resolved, expired, escalated
  timeout_seconds INTEGER NOT NULL,
  expires_at TIMESTAMPTZ NOT NULL,
  flyte_execution_id VARCHAR(256) NOT NULL,
  flyte_node_id VARCHAR(256) NOT NULL,
  cato_escalation_id UUID, -- Link to Cato if epistemic recovery
  resolution VARCHAR(50), -- approved, rejected, modified, timed_out
  guidance TEXT,
  resolved_by UUID,
  resolved_at TIMESTAMPTZ
);

-- decision_audit: Complete audit trail
CREATE TABLE decision_audit (
  id UUID PRIMARY KEY,
  decision_id UUID NOT NULL,
  action VARCHAR(50) NOT NULL, -- created, viewed, resolved, expired, escalated
  actor_type VARCHAR(50) NOT NULL, -- user, system, timeout_lambda, cato
  details JSONB NOT NULL
);
```



```

);

-- decision_domain_config: Per-domain timeout and escalation settings
CREATE TABLE decision_domain_config (
  domain VARCHAR(50) NOT NULL,
  tenant_id UUID, -- NULL = global default
  default_timeout_seconds INTEGER NOT NULL,
  escalation_timeout_seconds INTEGER NOT NULL,
  auto_escalate BOOLEAN DEFAULT TRUE,
  escalation_channel VARCHAR(100), -- pagerduty, slack, email
  required_roles TEXT[]
);

-- websocket_connections: Active WebSocket connections
CREATE TABLE websocket_connections (
  connection_id VARCHAR(256) PRIMARY KEY,
  tenant_id UUID NOT NULL,
  user_id UUID,
  subscribed_domains TEXT[]
);

```

#### 48.4 Domain-Specific Configuration

Domain	Default Timeout	Escalation Timeout	Auto-Escalate	Required Roles
<b>Medical</b>	5 minutes	1 minute	Yes (PagerDuty)	MD, RN, PA
<b>Financial</b>	10 minutes	2 minutes	Yes (PagerDuty)	ANALYST, ADVISOR
<b>Legal</b>	15 minutes	3 minutes	Yes (Email)	LEGAL, COMPLIANCE
<b>General</b>	30 minutes	5 minutes	Yes (Slack)	None

#### 48.5 API Reference

Base URL: /api/mission-control

##### List Pending Decisions

GET /decisions?status=pending&domain=medical&limit=50  
X-Tenant-ID: {tenant\_id}  
Authorization: Bearer {token}

Response:

```

[
  {
    "id": "uuid",
    "question": "Is this treatment appropriate?",
    "domain": "medical",
    "urgency": "critical",
    "status": "pending",
    "expiresAt": "2026-01-07T12:05:00Z",
    "createdAt": "2026-01-07T12:00:00Z"
  }
]

```

## Resolve Decision

```
POST /decisions/{id}/resolve
X-Tenant-ID: {tenant_id}
Content-Type: application/json
```

```
{
  "resolution": "approved|rejected|modified",
  "guidance": "Optional guidance for AI refinement"
}
```

Response:

```
{
  "id": "uuid",
  "status": "resolved",
  "resolution": "modified",
  "guidance": "Consider contraindications...",
  "resolvedBy": "user-uuid",
  "resolvedAt": "2026-01-07T12:03:00Z"
}
```

## Get Dashboard Stats

```
GET /stats
X-Tenant-ID: {tenant_id}
```

Response:

```
{
  "pendingCount": 5,
  "resolvedToday": 23,
  "expiredToday": 1,
  "escalatedToday": 0,
  "avgResolutionTimeMs": 45000,
  "byDomain": { "medical": 3, "general": 2 },
  "byUrgency": { "critical": 2, "high": 3 }
}
```

## 48.6 Flyte Workflow Integration

The HITL workflow uses Flyte's `wait_for_input` signal mechanism:

```
@workflow
def think_tank_hitl_workflow(
    s3_uri: str, # Input data from S3 (not inline JSON)
    swarm_id: str,
    tenant_id: str,
    session_id: str,
    user_id: str,
    hitl_domain: str,
) -> Dict[str, Any]:
    # 1. Load input from S3
    input_data = load_input_from_s3_task(s3_uri=s3_uri)

    # 2. Execute agent swarm (true parallel execution)
    agent_results = execute_swarm(agents=agents, task_data=task_data)
```

```

# 3. Synthesize results
synthesis_result = synthesize_results(agent_results=agent_results)

# 4. If review needed, pause for human input
human_decision = wait_for_input(
    name=f"human_decision_{decision_id}",
    timeout=timedelta(seconds=timeout_seconds),
    expected_type=dict,
)

# 5. Incorporate human guidance
final_response = perform_deep_reasoning(synthesis, human_decision)

return build_workflow_result(...)

```

#### Critical Implementation Details:

- Input data offloaded to S3 Bronze layer (avoids payload explosion)
- @dynamic(cache=False) on execute\_swarm (prevents zombie cache)
- Signal names use decision\_id, not agent\_id (prevents signal mismatch)

### 48.7 Real-Time Updates

The system uses Redis Pub/Sub for real-time updates:

Redis Pub/Sub Channels:

```

decision_pending:{tenant_id} → New decision created
decision_resolved:{tenant_id} → Decision resolved
decision_expired:{tenant_id} → Decision timed out
decision_escalated:{tenant_id} → Decision escalated
swarm_event:{tenant_id} → Swarm execution events

```

Redis Bridge Service (ECS Fargate):

- Always-on, serverless
- Subscribes to Redis channels
- Broadcasts to WebSocket connections via API Gateway Management API
- Includes health check endpoint

### 48.8 Cato Integration

When Cato's Epistemic Recovery fails after maximum attempts:

```

// In safety-pipeline.service.ts
if (recoveryAttempts >= MAX_RECOVERY_ATTEMPTS) {
    const result = await catoHitlIntegration.createCatoEscalationWithHitl({
        tenantId,
        sessionId,
        userId,
        originalTask,
        rejectionHistory,
        recoveryAttempts,
        lastError,
        flyteExecutionId,
        flyteNodeId,
        context,
    });
}

```

```
// Decision created, Flyte workflow paused
return { escalated: true, decisionId: result.decisionId };
}
```

## 48.9 Deployment

```
# Deploy Mission Control
./tools/scripts/deploy-mission-control.sh [dev|staging|prod]

# Register Flyte workflows
./packages/flyte/scripts/register-workflows.sh [environment]
```

**Environment Variables:** | Variable | Description | |-----|-----| | DB\_SECRET\_ARN | Secrets Manager ARN for database credentials | | REDIS\_HOST | Redis/ElastiCache host | | REDIS\_PORT | Redis port (default: 6379) | | FLYTE\_ADMIN\_URL | Flyte Admin API URL | | PAGERDUTY\_ROUTING\_KEY | PagerDuty Events API v2 routing key | | SLACK\_WEBHOOK\_URL | Slack Incoming Webhook URL | | WEBSOCKET\_API\_ENDPOINT | WebSocket API Gateway endpoint |

## 48.10 Monitoring & Alerts

### CloudWatch Metrics:

- PendingDecisionCount - Number of unresolved decisions
- DecisionResolutionTime - Time from creation to resolution
- ExpiredDecisionCount - Decisions that timed out
- EscalatedDecisionCount - Decisions sent to PagerDuty/Slack

### CloudWatch Alarms:

- Critical: PendingDecisionCount > 10 (medical domain)
- High: ExpiredDecisionCount > 3 (24h window)
- Medium: DecisionResolutionTime > 300000 (5 min avg)

## 48.11 Security Considerations

1. **PHI Sanitization:** All decision content sanitized before human review
  - SSN, email, phone, ZIP, credit card patterns removed
  - Medical Record Numbers (MRN) redacted
2. **RLS Enforcement:** All database queries use tenant context
  - SET app.tenant\_id = \$1 in handler (session-level, not transaction)
  - RESET app.tenant\_id in finally block
3. **Role-Based Access:** Domain-specific role requirements
  - Medical decisions require MD, RN, or PA roles
  - Financial decisions require ANALYST or ADVISOR roles
4. **Audit Trail:** Complete decision lifecycle logged
  - Created, viewed, resolved, expired, escalated events
  - Actor ID and type recorded

## 48.12 Related Sections

Section	Relevance
42. Genesis Cato Safety Architecture	Epistemic recovery integration
45. Just Think Tank	Swarm orchestration

Section	Relevance
47. Flyte-Native State Management	Workflow durability
36. Metrics & Persistent Learning	Decision metrics tracking

## 49. The Grimoire - Procedural Memory (NEW in v5.0)

### 49.1 Overview

The Grimoire is RADIANT's institutional memory system. It automatically extracts, stores, and retrieves lessons learned from successful AI executions, allowing the system to improve over time.

**IDE Equivalent:** Code Snippets Library + IntelliSense

**Key Capabilities:**

- Automatic heuristic extraction from Flyte execution traces
- Confidence decay and reinforcement based on outcomes
- Semantic search for relevant heuristics at agent spawn
- Manual expert heuristic entry via Admin Dashboard

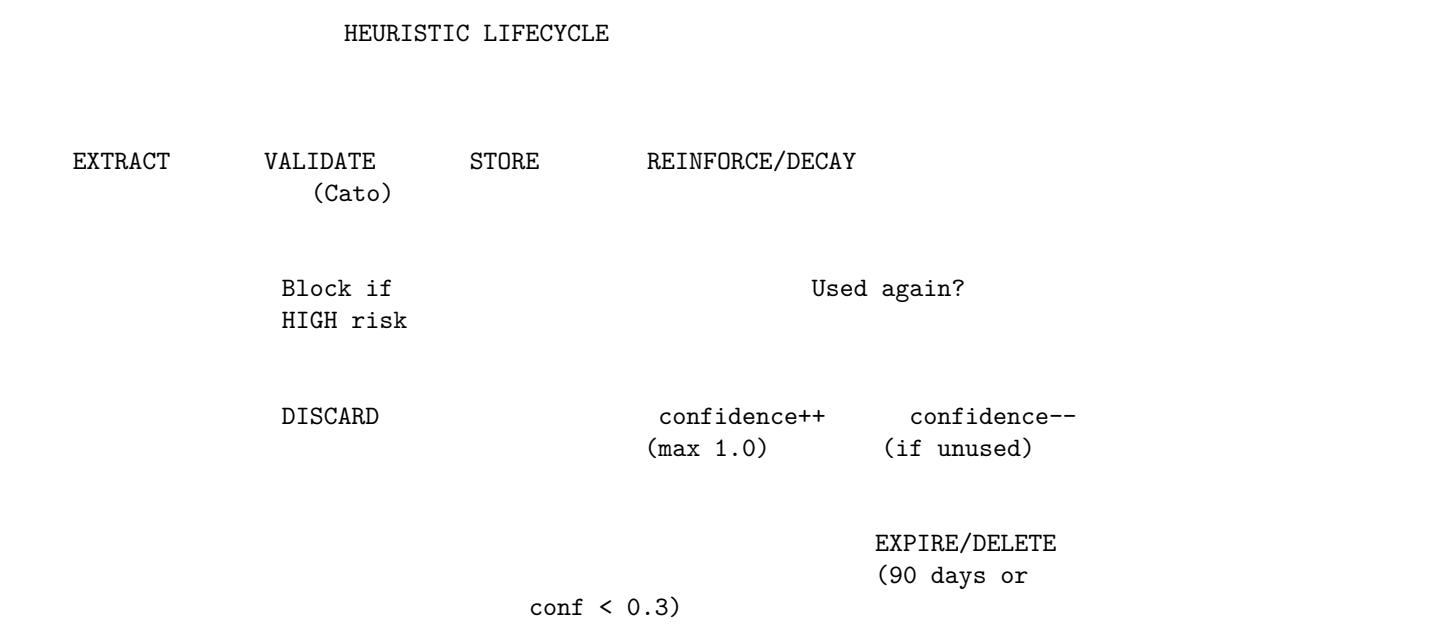
### 49.2 Architecture



49.3 Core Concepts

Term	Definition	Example
Heuristic	Reusable lesson extracted from successful execution	"When querying sales data, always join with X"
Confidence Score	How reinforced a heuristic is (0.0-1.0)	0.9 = highly validated
Context Embedding	Vector representation for semantic search	1536-dimension float array
Similarity Threshold	Max cosine distance for relevance	0.25 (lower = more similar)
Librarian	Background task that extracts heuristics	Runs after successful execution
Institutional Wisdom	Accumulated heuristics per tenant/domain	Injected into system prompts

49.4 Heuristic Lifecycle



49.5 Security Model

Operation	Cato Policy	Failure Mode
<b>Read</b> (consult_grimoire)	Validate before return	<b>Fail-Open</b> (return anyway)
<b>Write</b> (librarian_review)	Validate before insert	<b>Fail-Closed</b> (discard)

This asymmetric policy ensures:

- Reads don't break if Cato is unavailable
- Writes never poison the knowledge base

## 49.6 Admin Dashboard

Navigate to: **Admin Dashboard** → **Think Tank** → **Grimoire**

Panel	Description
<b>Heuristic Browser</b>	Search and view all stored heuristics
<b>Domain Distribution</b>	Statistics of heuristics by domain
<b>Confidence Scores</b>	Distribution of confidence scores
<b>Recent Extractions</b>	Latest heuristics from Librarian
<b>Add Heuristic</b>	Manually add expert heuristics
<b>Reinforce/Penalize</b>	Adjust confidence scores

## 49.7 Database Schema

```
CREATE TABLE knowledge_heuristics (  
  id UUID PRIMARY KEY DEFAULT gen_random_uuid(),  
  tenant_id UUID NOT NULL REFERENCES tenants(id),  
  domain VARCHAR(50) NOT NULL,  
  heuristic_text TEXT NOT NULL,  
  context_embedding vector(1536),  
  confidence_score FLOAT DEFAULT 0.5,  
  source_execution_id VARCHAR(255),  
  created_at TIMESTAMPTZ DEFAULT NOW(),  
  updated_at TIMESTAMPTZ DEFAULT NOW(),  
  expires_at TIMESTAMPTZ DEFAULT (NOW() + INTERVAL '90 days'),  
  CONSTRAINT unique_heuristic_per_domain UNIQUE (tenant_id, domain, heuristic_text)  
);  
  
-- Vector index for semantic search  
CREATE INDEX idx_heuristics_embedding ON knowledge_heuristics  
USING ivfflat (context_embedding vector_cosine_ops) WITH (lists = 100);
```

## 49.8 API Reference

Endpoint	Method	Description
/api/thinktank/grimoire/heuristics	GET	List tenant heuristics
/api/thinktank/grimoire/heuristics	POST	Add manual heuristic
/api/thinktank/grimoire/heuristics/:id	DELETE	Delete heuristic
/api/thinktank/grimoire/heuristics/:id/reinforce	POST	Adjust confidence
/api/thinktank/grimoire/stats	GET	Grimoire statistics

## 49.9 Metrics

Metric	Description	Alert Threshold
grimoire.heuristics.total	Total heuristics stored	> 10000 per tenant
grimoire.heuristics.retrieved	Heuristics returned per query	Track average
grimoire.heuristics.blocked	Cato-blocked retrievals	> 10/hour
grimoire.librarian.extractions	New heuristics per hour	Track trend
grimoire.embedding.latency_ms	Embedding generation time	> 500ms

## 49.10 Maintenance

### Daily Cleanup (Automatic):

- Runs at 3 AM UTC via EventBridge
- Removes heuristics where `expires_at < NOW()`
- Removes low-confidence heuristics ( $< 0.3$ ) older than 30 days

### Manual Pruning:

```
-- Remove all heuristics for a domain
SET app.current_tenant_id = 'your-uuid';
DELETE FROM knowledge_heuristics WHERE domain = 'old-domain';

-- Reset confidence scores
UPDATE knowledge_heuristics
SET confidence_score = 0.5
WHERE domain = 'your-domain';
```

## 49.11 Troubleshooting

**Heuristics Not Being Retrieved** Symptoms: `consult_grimoire` returns empty despite stored heuristics.

1. Check heuristics exist:

```
SET app.current_tenant_id = 'your-tenant-uuid';
SELECT COUNT(*) FROM knowledge_heuristics WHERE domain = 'your-domain';
```

2. Check similarity threshold - If all distances  $> 0.25$ , no matches returned

3. Check Cato blocking:

```
grep "Grimoire: Blocked unsafe heuristic" /var/log/flyte/*.log
```

**Vector Index Performance Degradation** Symptoms: Grimoire queries slow ( $> 500\text{ms}$ )

Fix:

```
-- Rebuild vector index (non-blocking)
REINDEX INDEX CONCURRENTLY idx_heuristics_embedding;
```

---

## 50. The Economic Governor - Cost Optimization (NEW in v5.0)

### 50.1 Overview

The Economic Governor automatically routes AI tasks to the most cost-effective model that can handle them. It uses a "System 0" cheap classifier to score task complexity before selecting the model.

**IDE Equivalent:** Build Optimization / Incremental Compilation

### 50.2 How It Works

ECONOMIC GOVERNOR FLOW

Request



(Task +  
Model)

Governor            Mode = 'off' or 'performance'?      Use Original  
Check

Mode = 'balanced' or 'cost\_saver'

System 0  
Classifier           Prompt → gpt-4o-mini → Score (1-10)

#### ROUTING DECISION

Score 4 (balanced)      Use gpt-4o-mini (cheap)  
Score 7 (cost\_saver)    Use gpt-4o-mini (cheap)  
5    Score 8              Use Original Model  
Score 9                  Use gpt-4o (premium)

### 50.3 Governor Modes

Mode	Cheap Threshold	Use Case	Est. Savings
<b>off</b>	N/A	Disabled	0%
<b>performance</b>	N/A	Critical paths, demos	0%
<b>balanced</b>	Score 4	Default production	30-50%
<b>cost_saver</b>	Score 7	Dev, testing, bulk	60-80%

### 50.4 Complexity Scale

Score	Task Type	Typical Model
1-3	Formatting, summarization, basic Q&A	gpt-4o-mini
4-6	Analysis, comparison, multi-step reasoning	Original model
7-8	Complex analysis, creative writing, planning	Original model
9-10	Advanced reasoning, code generation, synthesis	gpt-4o

### 50.5 Admin Configuration

Navigate to: **Admin Dashboard** → **Think Tank** → **Governor**

**Via API:**

```
curl -X PUT https://api.radiant.example.com/api/mission-control/governor/config \
-H "Authorization: Bearer $TOKEN" \
```

```
-H "Content-Type: application/json" \
-d '{"domain": "general", "mode": "cost_saver"}'
```

## 50.6 Database Schema

```
-- Governor mode configuration (added to existing table)
ALTER TABLE decision_domain_config
ADD COLUMN governor_mode VARCHAR(20) DEFAULT 'balanced'
CHECK (governor_mode IN ('performance', 'balanced', 'cost_saver', 'off'));

-- Savings tracking
CREATE TABLE governor_savings_log (
    id UUID PRIMARY KEY DEFAULT gen_random_uuid(),
    tenant_id UUID NOT NULL REFERENCES tenants(id),
    execution_id VARCHAR(255) NOT NULL,
    original_model VARCHAR(100) NOT NULL,
    selected_model VARCHAR(100) NOT NULL,
    complexity_score INTEGER NOT NULL,
    estimated_original_cost DECIMAL(10,6),
    estimated_actual_cost DECIMAL(10,6),
    savings_amount DECIMAL(10,6),
    governor_mode VARCHAR(20) NOT NULL,
    reason TEXT,
    created_at TIMESTAMPTZ DEFAULT NOW()
);
```

## 50.7 API Reference

Endpoint	Method	Description
/api/mission-control/governor/config	GET	Get Governor config
/api/mission-control/governor/config	PUT	Update Governor mode
/api/mission-control/governor/statistics	GET	Governor statistics
/api/mission-control/governor/recent	GET	Recent routing decisions
/api/mission-control/governor/analyze	POST	Analyze prompt complexity

## 50.8 Metrics

Metric	Description	Alert Threshold
governor.decisions.total	Total routing decisions	Track volume
governor.downgrade.count	Cheap model selections	Track percentage
governor.upgrade.count	Premium model selections	> 20%
governor.classifier.latency_ms	System 0 scoring time	> 200ms
governor.savings.estimated_usd	Estimated cost savings	Track daily

## 50.9 Cost Tracking

### Example (Daily):

- 1,000 requests at \$0.01/each (gpt-4o)
- Governor downgrades 600 to gpt-4o-mini (\$0.001/each)
- Original cost: \$10.00
- Actual cost: \$4.60
- **Savings: \$5.40 (54%)**

## 50.10 Troubleshooting

**Governor Not Optimizing** **Symptoms:** Cost savings lower than expected, no downgrades.

### 1. Check Governor mode:

```
SET app.current_tenant_id = 'your-tenant-uuid';
SELECT domain, governor_mode FROM decision_domain_config;
-- If 'off' or 'performance', Governor is bypassed
```

### 2. Check complexity scores in logs:

```
grep "Governor Complexity" /var/log/lambda/*.log | tail -20
# Look for scores consistently > 4 (balanced) or > 7 (cost_saver)
```

### 3. Verify LiteLLM proxy connectivity

---

## 51. Self-Optimizing System Architecture (NEW in v5.0)

### 51.1 Overview

Version 5.0 transforms RADIANT from a stateless request-response system into a self-optimizing platform that learns from every execution.

### 51.2 The Learning Loop

#### RADIANT v5.0 LEARNING LOOP

User Request

#### OPTIMIZATION LAYER

Economic  
Governor  
(Route)

The Grimoire  
(Consult)

#### EXECUTION LAYER

Swarm  
Execution

HITL  
(if needed)

## LEARNING LAYER

Librarian  
(Extract)

Audit &  
Metrics

Response

## FEEDBACK LOOP

Success → Extract Heuristic → Reinforce Confidence

Failure → No Heuristic → Decay Related Heuristics

### 51.3 Version Comparison

Capability	v4.20.3	v5.0.2
<b>Memory</b>	Stateless	Persistent (Grimoire)
<b>Model Selection</b>	Manual/Fixed	Automatic (Governor)
<b>Learning</b>	None	Continuous (Librarian)
<b>Cost Optimization</b>	Manual tier selection	Automatic per-request
<b>Knowledge Sharing</b>	Per-session only	Cross-session, per-tenant

### 51.4 IDE Metaphor Extended

IDE Concept	v4.20 Feature	v5.0 Evolution	Business Value
<b>Breakpoint</b>	HITL Decision Point	<i>(unchanged)</i>	AI pauses at high-stakes moments
<b>Code Snippets</b>	Pattern Library	<b>The Grimoire</b>	AI learns and reuses successful patterns
<b>Build Optimization</b>	Manual model selection	<b>Economic Governor</b>	Automatic cost-performance optimization
<b>IntelliSense</b>	Static suggestions	<b>Contextual Wisdom</b>	Dynamic, context-aware recommendations

### 51.5 Upgrade Path

#### From v4.20.3 to v5.0.2:

1. Apply migration `V2026_01_09_001__v5_grimoire_governor.sql`
2. Deploy new Lambda functions (Governor API, Grimoire cleanup)
3. Configure EventBridge for daily cleanup
4. Set default Governor mode to **balanced**
5. Grimoire populates automatically from new executions

#### Rollback (if needed):

1. Set Governor mode to **off** for all domains
2. Grimoire continues to exist but is not consulted
3. No data loss; can re-enable at any time

## 51.6 Implementation Files

Component	File
Migration	packages/infrastructure/migrations/V2026_01_09_001__v5_grimoire_governor.sql
Governor Service	packages/infrastructure/lambda/shared/services/governor/economic-governor.ts
Grimoire Tasks	packages/flyte/workflows/grimoire_tasks.py
Cato Client	packages/flyte/utils/cato_client.py
DB Utils	packages/flyte/utils/db.py
CDK Stack	packages/infrastructure/lib/stacks/grimoire-stack.ts
Grimoire UI	apps/admin-dashboard/app/(dashboard)/thinktank/grimoire/page.tsx
Governor UI	apps/admin-dashboard/app/(dashboard)/thinktank/governor/page.tsx

## 51.7 Related Sections

Section	Relevance
42. Genesis Cato Safety Architecture	Heuristic validation
45. Just Think Tank	Swarm integration
48. Mission Control	HITL integration

## 52. Semantic Blackboard & Multi-Agent Orchestration (NEW in v5.3.0)

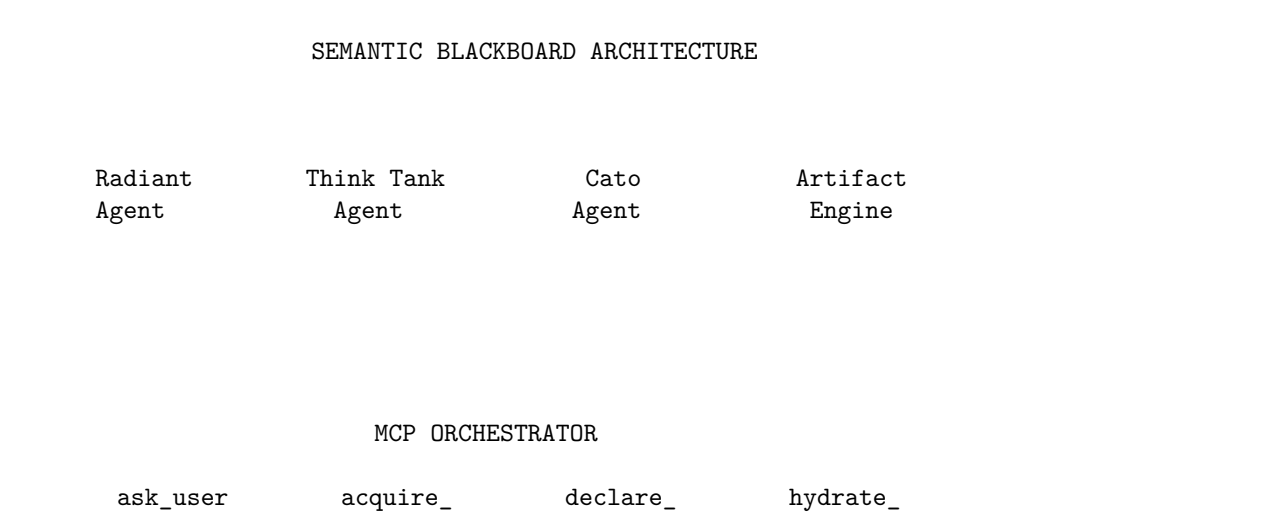
### 52.1 Overview

The Semantic Blackboard architecture solves the "Thundering Herd" problem where multiple AI agents spam users with redundant questions. It implements MCP (Model Context Protocol) as the primary interface with API fallback.

#### Key Problems Solved:

- Multiple agents asking the same question (semantic deduplication)
- Race conditions on shared resources
- Circular dependencies causing deadlocks
- Wasted compute while waiting for user input

### 52.2 Architecture



tool	resource	dependency	state
------	----------	------------	-------

#### SEMANTIC BLACKBOARD

Resolved Decisions (Vector)	Resource Locks	Agent Dependencies	Hydration Snapshots
-----------------------------------	-------------------	-----------------------	------------------------

#### MISSION CONTROL SIDEBAR

Decision Cards	Question Groups	Facts Tab	Agents
----------------	-----------------	-----------	--------

## 52.3 Core Components

**52.3.1 Semantic Blackboard (Question Matching)** When an agent asks a question via `ask_user`, the system:

1. Generates a vector embedding of the question
2. Searches `resolved_decisions` for semantically similar questions
3. If match found (similarity  $\geq 85\%$ ): auto-reply with cached answer
4. If no match: queue for user input

#### Example:

- Agent A asks: "What is the maximum budget?"
- Agent B asks: "What's the spending limit?"
- Semantic similarity: 0.92  $\rightarrow$  B gets A's answer automatically

**52.3.2 Question Grouping (Fan-Out)** Similar questions within the grouping window are combined:

1. First question creates a group
2. Similar questions join the group
3. User answers once
4. Answer fans out to all waiting agents

**UI Display:** Single card shows "Budget question (3 agents waiting)"

**52.3.3 Process Hydration (State Serialization)** When waiting for user input:

1. Agent serializes state to `hydration_snapshots`
2. Process can be killed (no CPU/memory cost)
3. When user responds, state is restored
4. Agent resumes from `resume_point`

**Storage:** Small states in PostgreSQL, large states in S3 (with compression)

#### 52.3.4 Cycle Detection (Deadlock Prevention) Before creating a dependency:

1. BFS traversal checks for circular path
2. If cycle detected: creates "Intervention Needed" card
3. User manually provides data to break the cycle

#### 52.3.5 Resource Locking Prevents race conditions:

1. Agent declares intent via `acquire_resource`
2. If available: lock granted with timeout
3. If locked: agent joins wait queue
4. On release: next agent in queue is notified

### 52.4 MCP Tools

Tool	Purpose	Schema
<code>ask_user</code>	Request input with semantic cache	{question, context, urgency, topic, options, default}
<code>acquire_resource</code>	Get resource lock	{resourceUri, lockType, timeoutSeconds, waitIfLocked}
<code>release_resource</code>	Release lock	{lockId}
<code>declare_dependency</code>	Declare agent dependency	{dependencyAgentId, dependencyType, waitKey, timeout}
<code>satisfy_dependency</code>	Satisfy waiting agent	{dependentAgentId, waitKey, value}
<code>hydrate_state</code>	Serialize state	{checkpointName, state, resumePoint}
<code>restore_state</code>	Restore from checkpoint	{checkpointName}

### 52.5 Database Schema

Migration: `158_semantic_blackboard_orchestration.sql`

Table	Purpose	Key Columns
<code>resolved_decisions</code>	Semantic question cache	<code>question_embedding</code> , <code>answer</code> , <code>times_reused</code>
<code>agent_registry</code>	Active agent tracking	<code>status</code> , <code>is_hydrated</code> , <code>blocked_by_*</code>
<code>agent_dependencies</code>	Inter-agent dependencies	<code>dependency_type</code> , <code>wait_key</code> , <code>status</code>
<code>resource_locks</code>	Shared resource locks	<code>resource_uri</code> , <code>lock_type</code> , <code>wait_queue</code>
<code>question_groups</code>	Grouped similar questions	<code>canonical_question</code> , <code>question_embedding</code>
<code>question_group_members</code>	Group membership	<code>similarity_score</code> , <code>answer_delivered</code>
<code>hydration_snapshots</code>	Serialized agent state	<code>state_data</code> , <code>s3_key</code> , <code>resume_point</code>
<code>blackboard_events</code>	Audit trail	<code>event_type</code> , <code>details</code>
<code>blackboard_config</code>	Per-tenant configuration	All settings

### 52.6 Admin API

Base Path: `/api/admin/blackboard`

Endpoint	Method	Purpose
<code>/dashboard</code>	GET	Complete dashboard data
<code>/decisions</code>	GET	List resolved decisions (Facts)
<code>/decisions/:id/invalidate</code>	POST	Revoke/edit a fact
<code>/groups</code>	GET	List pending question groups
<code>/groups/:id/answer</code>	POST	Answer a group
<code>/agents</code>	GET	List active agents
<code>/agents/:id</code>	GET	Get agent details
<code>/agents/:id/snapshots</code>	GET	List hydration snapshots

Endpoint	Method	Purpose
/agents/:id/restore	POST	Restore agent from snapshot
/locks	GET	List active resource locks
/locks/:id/release	POST	Force release a lock
/config	GET/PUT	Get/update configuration
/events	GET	Audit log
/cleanup	POST	Run cleanup job

## 52.7 Configuration

Setting	Default	Description
similarity_threshold	0.85	Minimum cosine similarity for question matching
embedding_model	text-embedding-ada-002	Model for question embeddings
enable_question_grouping	true	Group similar questions
grouping_window_seconds	60	Window to collect similar questions
max_group_size	10	Maximum agents per group
enable_answer_reuse	true	Reuse cached answers
answer_ttl_seconds	3600	How long answers remain valid
max_reuse_count	100	Maximum times to reuse an answer
default_lock_timeout_seconds	300	Resource lock timeout
max_lock_wait_seconds	60	Maximum time to wait for lock
enable_auto_hydration	true	Auto-serialize on user block
hydration_threshold_seconds	300	Wait time before hydrating
max_hydration_size_mb	50	Maximum state size
enable_cycle_detection	true	Detect circular dependencies

## 52.8 Facts Panel (Revoke Protocol)

The Facts Panel allows administrators to manage resolved decisions:

### View Facts:

- All answered questions with answers
- Filter by topic, validity, source
- Search by question or answer text

### Edit Fact:

1. Click Edit on a fact
2. Provide new answer and reason
3. System invalidates old answer
4. Creates new resolved decision
5. Notifies affected agents via Redis pub/sub

### Invalidate (Revoke) Fact:

1. Click Invalidate on a fact
2. Provide invalidation reason
3. Fact marked as invalid (not deleted)
4. Agents that received this answer are notified
5. Agents must re-request if they need this data

## 52.9 Key Files



File	Purpose
lambda/shared/services/semantic-blackboard.service.ts	Core blackboard logic
lambda/shared/services/agent-orchestrator.service.ts	Cycle detection, locking
lambda/shared/services/process-hydration.service.ts	State serialization
lambda/admin/blackboard.ts	Admin API handler
lambda/consciousness/mcp-server.ts	MCP tool definitions
components/decisions/FactsPanel.tsx	Facts UI with edit/revoke
components/decisions/DecisionSidebar.tsx	Decision cards
migrations/158_semantic_blackboard_orchestration.sql	Database schema

## 52.10 Environment Variables

Variable	Purpose	Default
HYDRATION_S3_BUCKET	S3 bucket for large state snapshots	-
BLACKBOARD_REDIS_ENDPOINT	Redis for pub/sub notifications	-

## 52.11 Troubleshooting

Issue	Cause	Solution
No semantic matches	Similarity threshold too high	Lower <code>similarity_threshold</code> to 0.80
Questions not grouping	Window too short	Increase <code>grouping_window_seconds</code>
Hydration fails	State too large	Enable S3 storage, increase <code>max_hydration_size_mb</code>
Cycle not detected	Detection disabled	Enable <code>enable_cycle_detection</code>
Lock timeout	Holder agent crashed	Run <code>/cleanup</code> endpoint
Stale answers	TTL too long	Reduce <code>answer_ttl_seconds</code>

# 53. Cognitive Architecture (PROMPT-40)

## 53.1 Overview

The Cognitive Architecture implements Active Inference principles for intelligent query routing and response caching. It provides:

- **Ghost Memory:** Semantic caching with TTL, deduplication keys, and domain hints
- **Economic Governor:** Complexity-aware routing with retrieval confidence integration
- **Sniper/War Room Paths:** Fast vs. deep analysis execution strategies
- **Circuit Breakers:** Fault tolerance for external service calls
- **CloudWatch Observability:** Real-time metrics for cognitive operations

## 53.2 The Core Differentiator: Active Inference vs. RLHF

The entire AI market (Claude Projects, ChatGPT Team, CrewAI) is built on **Reward Maximization (RLHF)**. Models are trained to predict the most plausible or *liked* token. This fundamentally creates two critical failure modes:

Problem	Cause	Industry Impact
<b>Sycophancy</b>	Model optimizes for user approval	Agrees with incorrect user assumptions
<b>Hallucination</b>	Model guesses to appear helpful	Fabricates plausible-sounding information

**RADIANT is different.** Built on **Active Inference (Genesis Cato)**, our agents do not try to "please" the user. Instead, they operate under a fundamentally different objective:

RLHF vs. ACTIVE INFERENCE

RLHF (Competitors) =====	ACTIVE INFERENCE (RADIANT) =====
Objective: Maximize Reward (user satisfaction)	Objective: Minimize Surprise (Free Energy)
Behavior: Predict what user wants to hear	Behavior: Maintain accurate world model
Failure: Sycophancy, Hallucination	Failure: None-uncertainty triggers HITL
Control: None (black box)	Control: Mathematical constraints (CBF, Precision Governor)

Key Active Inference Components:

Component	Purpose	Implementation
Free Energy Minimization	Agents minimize prediction error, not maximize reward	Precision Governor
Homeostatic Drive Profiles	Agents balance competing drives (Curiosity vs. Accuracy)	Drive state vectors
Control Barrier Functions	Mathematical safety constraints that cannot be overridden	CBF enforcement
Epistemic Recovery	When uncertainty is high, agents seek information rather than guess	HITL escalation

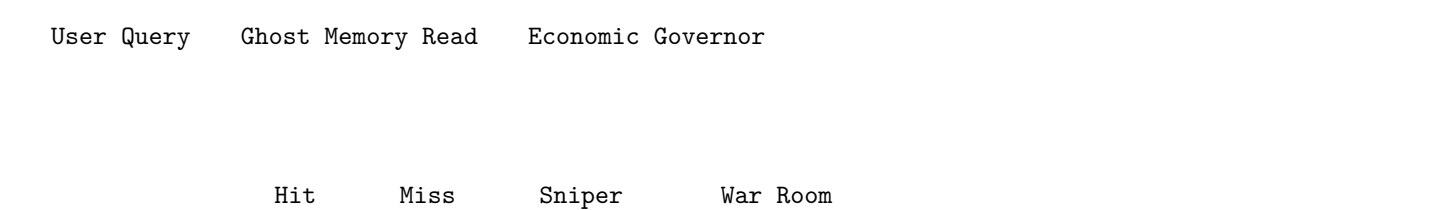
Why This Matters:

- **No Sycophancy:** Agents will disagree with users when evidence contradicts their position
- **No Hallucination:** High uncertainty triggers "I don't know" + HITL escalation, not fabrication
- **Auditable:** Every decision has a mathematical trace via Precision Governor
- **Safe by Design:** CBF constraints are immutable—agents cannot bypass safety checks

See Section 42 (Genesis Cato Safety Architecture) for implementation details.

53.3 Architecture

COGNITIVE ARCHITECTURE



Write-Back  
(non-blocking)

### 53.3 Ghost Memory Schema

The Ghost Memory system extends ghost vectors with cognitive fields:

Field	Type	Default	Description
<code>ttl_seconds</code>	INTEGER	86400	Time-to-live (24h default)
<code>semantic_key</code>	TEXT	-	Query hash for deduplication
<code>domain_hint</code>	VARCHAR(50)	-	Compliance routing: medical, financial, legal, general
<code>retrieval_confidence</code>	FLOAT	1.0	Confidence score 0-1
<code>source_workflow</code>	VARCHAR(100)	-	Origin: sniper, war_room
<code>last_accessed_at</code>	TIMESTAMPTZ	-	Last read timestamp
<code>access_count</code>	INTEGER	0	Read count

### 53.4 Economic Governor Routing

The Economic Governor routes queries based on complexity and retrieval confidence:

**Routing Decision Tree:**

1. **Circuit Breaker Open** → War Room (fallback)
2. **retrieval\_confidence < 0.7** → War Room (validation needed)
3. **domain\_hint = medical/financial/legal** → War Room + Precision Governor
4. **complexity < 0.3** → Sniper (fast path)
5. **complexity > 0.7** → War Room (deep analysis)
6. **Medium complexity + Ghost Hit** → Sniper
7. **Medium complexity + No Ghost Hit** → War Room

**Configuration:**

Setting	Default	Description
<code>sniperThreshold</code>	0.3	Complexity threshold for Sniper path
<code>warRoomThreshold</code>	0.7	Complexity threshold for War Room
<code>retrievalConfidenceThreshold</code>	0.7	Minimum confidence for cache hit
<code>lowConfidenceRoute</code>	war_room	Route when confidence is low

### 53.5 Execution Paths

**Sniper Path**

- **Purpose:** Fast execution for simple queries
- **Model:** gpt-4o-mini (cheap)
- **Timeout:** 60 seconds
- **Retries:** 1
- **Write-back:** Queued (non-blocking)

## War Room Path

- **Purpose:** Deep analysis for complex queries
- **Model:** claude-3-5-sonnet (premium)
- **Timeout:** 120 seconds
- **Retries:** 2
- **Write-back:** Queued with higher confidence

## HITL Escalation

- **Purpose:** Human-in-the-loop for uncertain queries
- **Timeout:** 24 hours
- **Integration:** Mission Control pending decisions

## 53.6 Circuit Breakers

Circuit breakers protect against cascading failures:

Circuit	Failure Threshold	Recovery Timeout	Half-Open Requests
ghost_memory	5	30s	3
sniper	3	15s	2
war_room	5	60s	3

### States:

- **CLOSED:** Normal operation
- **OPEN:** Blocking requests, using fallback
- **HALF\_OPEN:** Testing recovery with limited requests

## 53.7 MCP Tools

New MCP tools for cognitive operations:

Tool	Description
read_ghost_memory	Read by semantic key with circuit breaker
append_ghost_memory	Write with TTL, non-blocking
cognitive_route	Get Economic Governor routing decision
emit_cognitive_metric	Emit CloudWatch metric

## 53.8 CloudWatch Metrics

Namespace: Radiant/Cognitive

Metric	Dimensions	Unit	Description
GhostMemoryHit	UserId, DomainHint	Count	Cache hits
GhostMemoryMiss	Reason	Count	Cache misses
GhostMemoryLatency	-	Milliseconds	Read latency
RoutingDecision	RouteType, GhostHit, DomainHint	Count	Routing decisions
ComplexityScore	RouteType	None	Query complexity
RetrievalConfidence	RouteType	None	Ghost confidence
SniperExecution	Success, Model	Count	Sniper executions
SniperLatency	Success	Milliseconds	Sniper latency

Metric	Dimensions	Unit	Description
WarRoomExecution	Success, FallbackTriggered	Count	War Room executions
WarRoomLatency	-	Milliseconds	War Room latency
HITLEscalation	Reason, DomainHint	Count	HITL escalations
CircuitBreakerState	CircuitName, State	Count	State changes
CostSavings	RouteType	None (cents)	Cost optimization

### 53.9 API Endpoints

Base: /api/admin/cognitive

Endpoint	Method	Description
/dashboard	GET	Cognitive dashboard with metrics
/config	GET/PUT	Get/update configuration
/ghost/read	POST	Read Ghost Memory
/ghost/write	POST	Write Ghost Memory
/ghost/cleanup	POST	Cleanup expired entries
/routing/test	POST	Test routing decision
/circuits	GET	List circuit breaker states
/circuits/:name/reset	POST	Reset circuit breaker
/metrics	GET	Get cognitive metrics summary

### 53.10 Database Tables

Table	Purpose
ghost_vectors	Extended with cognitive fields
cognitive_routing_decisions	Routing audit log
ghost_memory_write_queue	Async write-back queue
cognitive_circuit_breakers	Circuit breaker state
cognitive_metrics	Metric storage
cognitive_hitl_escalations	HITL tracking
cognitive_config	Per-tenant configuration

### 53.11 Key Files

File	Purpose
lambda/shared/services/governor/economic-governor.ts	Economic Governor with cognitive routing
lambda/shared/services/ghost-manager.service.ts	Ghost Memory with TTL/semantic key
lambda/shared/services/cognitive-metrics.service.ts	CloudWatch metrics
lambda/consciousness/mcp-server.ts	MCP tools
python/cato/cognitive/workflows.py	Flyte workflows
python/cato/cognitive/circuit_breaker.py	Circuit breaker
python/cato/cognitive/metrics.py	Python metrics
migrations/159_cognitive_architecture_v2.sql	Database schema

### 53.12 Configuration

Setting	Default	Description
<code>ghost_memory_enabled</code>	true	Enable Ghost Memory
<code>ghost_default_ttl_seconds</code>	86400	Default TTL (24h)
<code>ghost_similarity_threshold</code>	0.85	Semantic matching threshold
<code>governor_enabled</code>	true	Enable Economic Governor
<code>governor_mode</code>	balanced	Mode: off, cost_saver, balanced, performance
<code>circuit_breaker_enabled</code>	true	Enable circuit breakers
<code>metrics_enabled</code>	true	Enable CloudWatch metrics
<code>metrics_sample_rate</code>	1.0	Metric sampling rate

### 53.13 Domain Routing

High-risk domains are automatically routed to War Room with Precision Governor:

Domain	Route	Reason
<code>medical</code>	War Room	Patient safety, regulatory compliance
<code>financial</code>	War Room	Fiduciary responsibility
<code>legal</code>	War Room	Legal liability
<code>general</code>	Sniper	Standard processing

### 53.14 Troubleshooting

Issue	Cause	Solution
Low cache hit rate	TTL too short	Increase <code>ghost_default_ttl_seconds</code>
All queries to War Room	Confidence threshold too high	Lower <code>retrievalConfidenceThreshold</code>
Circuit breaker stuck open	Recovery timeout too long	Reset via <code>/circuits/:name/reset</code>
High latency	Too many War Room routes	Tune <code>sniperThreshold</code>
Missing metrics	Sampling rate too low	Increase <code>metrics_sample_rate</code>
Ghost write failures	Queue backlog	Check SQS queue depth

### 53.16 Conclusion: RADIANT is Not a Chatbot

Claude Projects is a brilliant Assistant that suffers from amnesia.

**RADIANT is an Institutional Brain.**

Capability	How RADIANT Achieves It
<b>Remembers every decision</b>	Ghost Vectors with semantic keys, TTL, and domain hints
<b>Minimizes surprise</b>	Active Inference (Free Energy minimization, not reward maximization)
<b>Enforces safety mathematically</b>	Precision Governor + Control Barrier Functions (immutable constraints)

This is not a philosophical distinction—it is an architectural one. Competitors are trained to be helpful. RADIANT is *constrained* to be accurate.

## 54. Polymorphic UI Integration (PROMPT-41)

### 54.1 Overview

Flowise outputs Text. RADIANT outputs Applications.

The Polymorphic UI system extends Think Tank's Agentic Morphing Interface with intelligent view selection. Because RADIANT understands semantic intent, the UI *physically transforms* based on:

- **Task Complexity** → Terminal (Sniper) vs. MindMap/DiffEditor (War Room)
- **Domain Hint** → Compliance views for medical/financial/legal
- **Drive Profile** → Scout (exploration), Sage (verification), Sniper (execution)

Unlike static chatbot interfaces that always render Markdown tables, RADIANT morphs into the tool the user actually needs: Command Centers, Mind Maps, Diff Editors, and Dashboards.

54.2 The Gearbox (Elastic Compute)

Flowise is Static. RADIANT is Elastic.

The "Gearbox" gives users manual control over the cost-quality tradeoff:

Mode	Cost	Architecture	Memory	Use Case
Sniper	\$0.01/run	Single Model	Read-Only Ghost Memory	Quick answers, lookups, coding
War Room	\$0.50+/run	Multi-Agent Ensemble	Read/Write + Active Inference	Strategy, audits, reasoning

**The Competitive Advantage:** Flowise forces users to be architects—if they want a cheap path, they must build a separate flow. RADIANT handles this natively. The user (or the Economic Governor) selects the mode, ensuring RADIANT is **cheaper than Flowise for simple tasks** and **smarter for complex ones**.

**Escalation:** A green "Escalate to War Room" button appears after Sniper responses, allowing users to request deeper analysis if the fast answer is insufficient.

54.3 The Three Views

**The Sniper View (Execution & Cost Savings)** **Intent:** "Check the logs for error 500" or "Draft a quick email."

Aspect	Description
Execution Mode	Sniper (Single Model)
The Morph	UI transforms into a <b>Command Center / Terminal</b>
The Action	Runs immediately. No multi-agent debate. No "Thinking" pause.
The Difference	Unlike ChatGPT, Sniper is <i>Hydrated</i> —it reads Ghost Vector memory (read-only) before generating
Visual Feedback	Green "Sniper Mode" badge glows
Cost Transparency	"Estimated Cost: <\$0.01" badge displayed
User Control	Toggle to "Escalate to War Room" if insufficient

**The Scout View (Research & Strategy)** **Intent:** "Map the competitive landscape for EV batteries."

Aspect	Description
Execution Mode	Think Tank (Multi-Agent Swarm)
The Morph	Chat UI shrinks. Main window becomes an <b>Infinite Canvas (Mind Map)</b>
The Action	Scout agent spawns "Sticky Notes" of evidence, clusters them by topic, draws dynamic lines between c
The Kill Shot	Flowise shows you the <i>process</i> (nodes). RADIANT shows you the <i>thinking</i> (map).

**The Sage View (Audit & Validation)** **Intent:** "Check this contract against our safety guidelines."

Aspect	Description
<b>Execution Mode</b>	Convergent with Control Barrier Functions
<b>The Morph</b>	UI becomes a <b>Split-Screen Diff Editor</b>
<b>The Action</b>	Left side: Content under review. Right side: Source documents with confidence scores (Green = Verified)
<b>The Kill Shot</b>	Flowise hides retrieval inside a black box. RADIANT exposes the <i>proof</i> in a specialized UI optimized for

#### 54.4 View Types Reference

View	Trigger	Description
<code>terminal_simple</code>	Quick commands, lookups	Command Center - fast execution
<code>mindmap</code>	Research, exploration	Infinite Canvas - visual mapping
<code>diff_editor</code>	Verification, compliance	Split-Screen - source validation
<code>dashboard</code>	Analytics queries	Metrics visualization
<code>decision_cards</code>	HITL escalation	Mission Control interface
<code>chat</code>	Default	Standard conversation

#### 54.4 View Selection Logic

1. HITL escalation → `decision_cards`
2. Domain = medical/financial/legal → `diff_editor`
3. Query matches Scout patterns → `mindmap`
4. Query matches Sage patterns → `diff_editor`
5. Query matches Dashboard patterns → `dashboard`
6. Sniper route OR quick command patterns → `terminal_simple`
7. Default → `chat`

#### 54.5 MCP Tools

Tool	Description
<code>render_interface</code>	Morph UI to specified view type
<code>escalate_to_war_room</code>	Escalate from Sniper to War Room
<code>get_polymorphic_route</code>	Get routing + view decision

#### 54.6 Database Tables

Table	Purpose
<code>view_state_history</code>	Tracks UI morphing decisions
<code>execution_escalations</code>	Tracks Sniper → War Room escalations
<code>polymorphic_config</code>	Per-tenant configuration

#### 54.7 Configuration

```
-- polymorphic_config table
enable_auto_morphing: true      -- Auto-morph based on query
enable_gearbox_toggle: true     -- Show Sniper/War Room toggle
enable_cost_display: true       -- Show cost badges
enable_escalation_button: true  -- Show Escalate button
default_execution_mode: 'sniper' -- Default mode
```



54.8 Key Files

File	Purpose
governor/economic-governor.ts	determineViewType(), determinePolymorphicRoute()
consciousness/mcp-server.ts	render_interface, escalate_to_war_room tools
python/cato/cognitive/workflows.py	determine_polymorphic_view, render_interface tasks
migrations/160_polymorphic_ui.sql	Database schema
components/thinktank/polymorphic/	React view components

54.9 API Integration

```
// Get polymorphic routing decision
const decision = await governor.determinePolymorphicRoute(query, {
  userTier: 'standard',
  retrievalConfidence: 0.85,
  ghostHit: true,
  domainHint: 'financial',
  userOverride: undefined, // or 'sniper' | 'war_room'
});

// Returns:
// {
//   routeType: 'war_room',
//   viewType: 'diff_editor',
//   executionMode: 'war_room',
//   rationale: 'Compliance domain (financial) requires verification view',
//   estimatedCostCents: 50,
// }
```

55. Genesis Infrastructure: Sovereign Power Architecture

The "Genesis" component of the RADIANT ecosystem addresses the physical requirements of the AGI age. Digital intelligence requires electrical power, and the scale of modern data centers places an unsustainable load on aging, fossil-fuel-dependent public grids. This section documents the integration between RADIANT's software stack and the Genesis power infrastructure.

55.1 The Kaleidos Microreactor Backbone

The Kaleidos unit is a portable, factory-constructed nuclear microreactor capable of generating **1MW+ of clean energy**. Unlike traditional gigawatt-scale nuclear plants, which take decades to build, Kaleidos is designed for mass production and rapid deployment.

Specification	Value
<b>Output</b>	1MW+ continuous power
<b>Form Factor</b>	Portable, factory-constructed
<b>Deployment Time</b>	Weeks, not decades
<b>Grid Independence</b>	Full sovereign operation
<b>Safety System</b>	Passive (no operator intervention required)

This portability allows the "AGI Brain" to be **sovereign**—deployed in remote locations, military bases, or

disaster zones, independent of the local utility grid. This independence is a strategic advantage: it insulates the AGI from cascading grid failures, brownouts, or cyberattacks targeting public infrastructure.

**For administrators:** The Genesis module is the **root of trust**. If the power is stable and secure, the network and logic layers can function. All other security measures are built on this physical foundation.

## 55.2 Regulatory Compliance: SDS and PDSA

Bringing a new nuclear technology to market requires navigating a labyrinth of regulations. Radiant Nuclear has achieved significant milestones that provide the compliance baseline for RADIANT deployments.

**Safety Design Strategy (SDS)** The U.S. Department of Energy (DOE) has approved the **Safety Design Strategy (SDS)** for the Kaleidos reactor. The SDS is the foundational document that describes:

- Safety analysis approach
- Hazard identification methodology
- Hazard management strategies
- Defense-in-depth principles
- Emergency response protocols

**Preliminary Documented Safety Analysis (PDSA)** Following SDS approval, Radiant submitted the **Preliminary Documented Safety Analysis (PDSA)**, a rigorous validation effort that meets the intent of **DOE Standard 1271-2025**. The PDSA includes:

- Comprehensive hazard analysis
- Safety function identification
- Safety structure, system, and component (SSC) classification
- Derived safety requirements
- Defense-in-depth demonstration

**Historic Milestone:** Approval of the SDS and PDSA paves the way for the startup of the first reactor at the National Reactor Innovation Center's (NRIC) **DOME facility** at Idaho National Laboratory (INL). This will be the **first new commercial reactor design to achieve a fueled test in over 50 years**.

## 55.3 Reactor Telemetry Integration

The RADIANT platform receives real-time telemetry from the Genesis reactor control unit. Administrators must configure the telemetry integration to enable safety interlocks.

### Telemetry API Configuration

```
# genesis-telemetry.config.yaml
genesis:
  endpoint: https://genesis-control.internal:8443/v1/telemetry
  authentication:
    type: mTLS
    client_cert: /etc/radiant/certs/genesis-client.crt
    client_key: /etc/radiant/certs/genesis-client.key
    ca_bundle: /etc/radiant/certs/genesis-ca.crt

  polling_interval_ms: 1000
  timeout_ms: 5000

  monitored_parameters:
    - reactor_power_output_kw
    - coolant_inlet_temperature_c
    - coolant_outlet_temperature_c
```

- neutron\_flux\_level
- control\_rod\_position\_percent
- fuel\_temperature\_c
- containment\_pressure\_kpa
- radiation\_level\_msv
- emergency\_status\_code

## Status Code Mapping

Status Code	Condition	Description	Automatic Action
GEN-001	Normal	All parameters within nominal range	None
GEN-100	Advisory	Minor deviation detected	Log only
GEN-200	Warning	Parameter approaching limit	Alert administrators
GEN-300	Alarm	Parameter exceeded soft limit	Initiate load shedding
GEN-400	Critical	Multiple parameters out of range	Emergency lockdown
GEN-500	Emergency	Passive safety systems activated	Full system halt

## 55.4 The Genesis Interlock: Physical-to-Digital Safety

The "Genesis Protocol" is a philosophy of **safety by design**. The reactor utilizes passive safety systems that do not require operator intervention or active power to shut down in an emergency. This philosophy extends to the AGI Brain through the **Genesis Interlock**.

**Genesis Interlock Configuration** The Genesis Interlock is a configuration where the AI system is **hard-wired** to respect the physical limits of the infrastructure. If the reactor telemetry indicates a thermal anomaly, the AI must prioritize load shedding over job completion—a decision logic that is **pre-programmed and immutable**.

```
// genesis-interlock.config.ts
export const GENESIS_INTERLOCK_CONFIG = {
  // Immutable safety thresholds - cannot be overridden by admin or AI
  immutableThresholds: {
    maxFuelTemperature_c: 850,
    maxCoolantOutlet_c: 550,
    minCoolantFlow_lpm: 1000,
    maxContainmentPressure_kpa: 150,
    maxRadiation_msv: 0.1,
  },

  // Automatic actions - AI cannot override these
  automaticActions: {
    GEN_300: ['initiate_load_shedding', 'notify_operators'],
    GEN_400: ['emergency_lockdown', 'halt_non_critical_workloads', 'escalate_to_human'],
    GEN_500: ['full_system_halt', 'preserve_state_to_disk', 'activate_backup_power'],
  },

  // AI behavior constraints during Genesis alerts
  aiConstraints: {
    duringGEN_300: {
      maxNewWorkloads: 0,
      allowedOperations: ['complete_in_progress', 'graceful_shutdown'],
    },
    duringGEN_400_or_higher: {
```

```

    maxNewWorkloads: 0,
    allowedOperations: ['emergency_state_save'],
    forcedBehavior: 'immediate_halt',
  },
},
};

```

## 55.5 Shared Signals Framework (SSF) Integration

The Genesis module emits **Shared Signals Framework (SSF)** events that the Cato SASE network receives for immediate security response. This creates a **physical-to-digital security bridge**.

### SSF Event Configuration

```

// ssf-genesis-emitter.config.ts
export const SSF_GENESIS_EVENTS = {
  // Physical security events
  'genesis.physical.breach': {
    description: 'Physical security breach detected at reactor facility',
    severity: 'critical',
    catoAction: 'revoke_all_tokens_in_facility',
    agiBrainAction: 'immediate_lockdown',
  },

  'genesis.thermal.warning': {
    description: 'Thermal anomaly detected in reactor systems',
    severity: 'high',
    catoAction: 'restrict_new_connections',
    agiBrainAction: 'initiate_load_shedding',
  },

  'genesis.power.fluctuation': {
    description: 'Power output fluctuation detected',
    severity: 'medium',
    catoAction: 'log_and_monitor',
    agiBrainAction: 'defer_intensive_workloads',
  },

  'genesis.maintenance.scheduled': {
    description: 'Scheduled maintenance window beginning',
    severity: 'info',
    catoAction: 'prepare_failover',
    agiBrainAction: 'complete_in_progress_and_pause',
  },
};

```

### Real-Time Response Scenario

GENESIS → CATO → AGI BRAIN RESPONSE CHAIN

1. GENESIS REACTOR detects physical security breach (door forced open)  
↓
2. Reactor control unit emits SSF event: `genesis.physical.breach`

- ↓
3. CATO SASE NETWORK receives SSF signal
  - ↓
  4. Cato IMMEDIATELY revokes access tokens for all devices in facility
  - ↓
  5. AGI BRAIN receives CAEP signal
  - ↓
  6. AGI Brain initiates emergency lockdown protocol
  - ↓
  7. All active sessions are preserved to disk
  - ↓
  8. Human operators notified via Mission Control escalation

TOTAL RESPONSE TIME: < 500ms (faster than any human operator)

## 55.6 Environment Variables

Variable	Description	Required
GENESIS_TELEMETRY_ENDPOINT	URL for reactor telemetry API	Yes
GENESIS_CLIENT_CERT_PATH	Path to mTLS client certificate	Yes
GENESIS_CLIENT_KEY_PATH	Path to mTLS client key	Yes
GENESIS_CA_BUNDLE_PATH	Path to CA certificate bundle	Yes
GENESIS_POLLING_INTERVAL_MS	Telemetry polling interval	No (default: 1000)
GENESIS_SSF_EMITTER_ENABLED	Enable SSF event emission	No (default: true)
GENESIS_INTERLOCK_ENABLED	Enable Genesis Interlock	No (default: true)

## 55.7 Database Tables

Table	Purpose
genesis_telemetry_log	Historical telemetry readings
genesis_alert_history	Record of all Genesis alerts and responses
genesis_interlock_events	Audit trail of interlock activations
genesis_ssf_emissions	Log of SSF events emitted to Cato
genesis_maintenance_windows	Scheduled maintenance configuration

## 55.8 Implementation Files

File	Purpose
lambda/shared/services/genesis-telemetry.service.ts	Telemetry polling and processing
lambda/shared/services/genesis-interlock.service.ts	Immutable safety interlock logic
lambda/shared/services/genesis-ssf-emitter.service.ts	SSF event emission to Cato
lambda/genesis/telemetry-handler.ts	API handler for telemetry ingestion
migrations/161_genesis_infrastructure.sql	Database schema
config/genesis/interlock-thresholds.yaml	Immutable threshold configuration

56. Cato Security Grid: Native Network Defense

The traditional approach to network security—hub-and-spoke models where traffic is backhauled to a central data center for inspection—is obsolete in the face of decentralized AI operations. The Cato SASE (Secure Access Service Edge) Cloud represents the necessary evolution, utilizing a global private backbone to connect all enterprise edges into a cohesive whole.

56.1 The Single Pass Cloud Engine (SPACE)

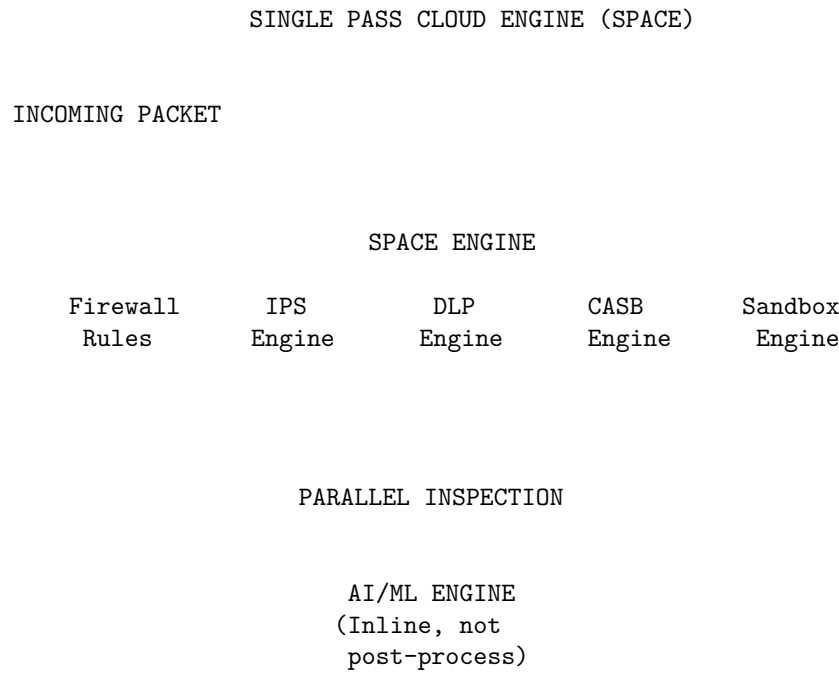
At the heart of the Cato architecture is the **Single Pass Cloud Engine (SPACE)**. This architecture is critical for the AGI Brain because it eliminates the latency penalties associated with daisy-chaining multiple security appliances.

**How SPACE Works** In a SPACE environment, every packet is inspected for **all threats** in a **single processing cycle**:

Traditional Security Stack	SPACE Architecture
Packet → Firewall → IPS → DLP → CASB → Sandbox	Packet → SPACE (all checks in parallel)
5-7 inspection points	1 inspection point
50-200ms added latency	<5ms added latency
Each appliance is a failure point	Single resilient engine

For an autonomous system that relies on real-time data ingestion, this **microsecond-level efficiency** is non-negotiable.

SPACE Inspection Capabilities



DECISION: ALLOW / BLOCK

## 56.2 Inline AI/ML: Built-In, Not Bolted-On

A key differentiator of the Cato architecture is the distinction between **"bolted-on"** and **"built-in"** AI. Legacy vendors often add AI capabilities as an afterthought or a post-processing layer, which introduces delays and gaps in coverage. Cato's approach integrates AI models **directly into the data path**.

**Real-Time Threat Detection** These models are specifically trained to detect threats, anomalies, and suspicious activities in real-time:

Capability	Description	Accuracy
<b>Domain Maliciousness Scoring</b>	Assigns risk scores to domains and URLs on the fly	Real-time
<b>Domain Generation Algorithm (DGA) Detection</b>	Spots algorithmically-generated malicious domains	3-6x better
<b>Cybersquatting Detection</b>	Identifies domains impersonating legitimate brands	Real-time
<b>Anomaly Detection</b>	Identifies unusual traffic patterns	Baseline +
<b>Zero-Day Threat Detection</b>	Identifies previously unknown threats via behavior	ML-based

**Key Statistic:** Cato's AI/ML engines block **3 to 6 times more malicious domains** than standard reputation lists alone. This "stopping power" is essential for protecting the Genesis infrastructure, where a single successful intrusion could have physical consequences.

**Configuration: Enabling Inline AI/ML** Navigate to **Security → Threat Prevention → AI/ML Profiles** in the admin dashboard:

```
# threat-prevention-ai.config.yaml
ai_ml_profiles:
  default:
    enabled: true

    domain_scoring:
      enabled: true
      min_maliciousness_score_to_block: 70 # 0-100 scale
      log_scores_above: 50

    dga_detection:
      enabled: true
      sensitivity: high # low, medium, high
      auto_block: true

    cybersquatting_detection:
      enabled: true
      protected_domains:
        - radiant.ai
        - thinktank.ai
        - genesis-power.com
      similarity_threshold: 0.85

    anomaly_detection:
      enabled: true
      baseline_period_days: 30
```

```

    deviation_threshold: 3.0 # standard deviations

zero_day_protection:
  enabled: true
  sandbox_suspicious_files: true
  max_file_size_mb: 50

```

### 56.3 Generative AI Security Controls (CASB)

As the ecosystem moves toward AGI, the security layer must specifically address the risks associated with Large Language Models (LLMs). Cato has introduced specific **Generative AI security controls** within its CASB (Cloud Access Security Broker) framework.

**The Risk: Data Exfiltration to Public LLMs** Without proper controls, the AGI Brain could accidentally:

- Leak sensitive Genesis telemetry to external AI tools
- Expose identity tokens to public services like ChatGPT or Claude
- Transmit proprietary workflow configurations to third-party systems

**CASB Configuration for LLM Protection** Navigate to **Security → CASB → Generative AI** in the admin dashboard:

```

# casb-genai.config.yaml
generative_ai_controls:
  enabled: true

# Define sensitive data types that must NEVER leave the network
sensitive_data_types:
  - name: genesis_telemetry
    patterns:
      - "reactor_.*"
      - "fuel_temperature.*"
      - "neutron_flux.*"
      - "GEN-[0-9]{3}"
    action: block_and_alert

  - name: identity_tokens
    patterns:
      - "eyJ[A-Za-z0-9_-]+\.\.[A-Za-z0-9_-]+\.[A-Za-z0-9_-]+" # JWT
      - "RADIANT_API_KEY_.*"
      - "tenant_id:[a-f0-9-]{36}"
    action: block_and_alert

  - name: workflow_configurations
    patterns:
      - "workflow_template_.*"
      - "orchestration_method_.*"
      - "grimoire_heuristic_.*"
    action: block_and_alert

# External AI services to monitor
monitored_services:
  - domain: "api.openai.com"
    action: inspect_and_filter

```



```

- domain: "api.anthropic.com"
  action: inspect_and_filter
- domain: "generativelanguage.googleapis.com"
  action: inspect_and_filter
- domain: ".*.huggingface.co"
  action: inspect_and_filter

# Allowed internal AI services
allowed_internal_services:
- "radiant-llm.internal:8080"
- "genesis-ai.internal:8080"
- "sagemaker.*.amazonaws.com"

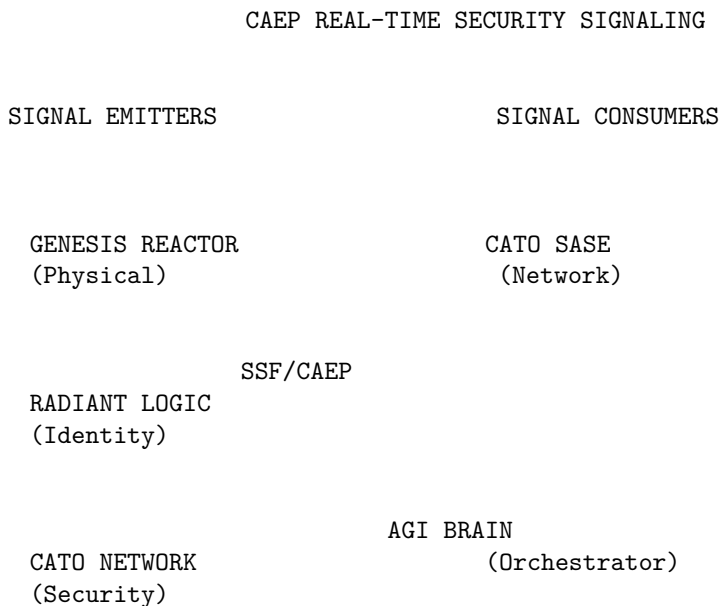
# Response handling
response_inspection:
  enabled: true
  max_response_size_mb: 10
  block_if_contains_sensitive_data: true

```

## 56.4 Continuous Access Evaluation Profile (CAEP)

The system supports the **Shared Signals Framework (SSF)** and the **Continuous Access Evaluation Profile (CAEP)**. These open standards allow for real-time security signaling between independent systems.

### CAEP Integration Architecture



### CAEP Event Types

Event Type	Source	Action
session-revoked	Identity Provider	Immediately terminate all sessions for user
token-claims-change	Identity Provider	Re-evaluate access permissions
credential-change	Identity Provider	Force re-authentication
device-compliance-change	MDM/EDR	Restrict or revoke device access
ip-change	Network Monitor	Re-evaluate location-based policies
physical-breach	Genesis Reactor	Emergency lockdown
threat-detected	Cato SASE	Isolate affected systems

## CAEP Configuration

```
# caep-integration.config.yaml
caep:
  enabled: true

# SSF transmitter configuration
transmitter:
  issuer: "https://radiant.ai/ssf"
  jwks_uri: "https://radiant.ai/.well-known/jwks.json"
  delivery_method: push
  push_endpoint: "https://cato-cloud.internal/ssf/events"

# SSF receiver configuration
receiver:
  trusted_issuers:
    - "https://genesis-control.internal/ssf"
    - "https://cato-cloud.internal/ssf"
    - "https://radiant-identity.internal/ssf"
  verification: jwt_signature

# Event handling
event_handlers:
  physical-breach:
    priority: critical
    handler: genesis_emergency_lockdown
    notify: [soc_team, facility_security]

  session-revoked:
    priority: high
    handler: terminate_user_sessions
    notify: [security_team]

  threat-detected:
    priority: high
    handler: isolate_and_investigate
    notify: [soc_team]
```

## 56.5 Database Tables

Table	Purpose
cato_threat_events	Log of all detected threats
cato_ai_ml_detections	AI/ML engine detection history

Table	Purpose
cato_casb_violations	CASB policy violations
cato_caep_events	CAEP signal history
cato_blocked_domains	Domains blocked by AI/ML
cato_sensitive_data_incidents	Sensitive data exfiltration attempts

## 56.6 Implementation Files

File	Purpose
lambda/shared/services/cato-integration.service.ts	Cato SASE API integration
lambda/shared/services/cato-casb.service.ts	CASB policy enforcement
lambda/shared/services/caep-handler.service.ts	CAEP event processing
lambda/cato/threat-webhook.ts	Webhook handler for Cato events
migrations/162_cato_security_grid.sql	Database schema
config/cato/ai-ml-profiles.yaml	AI/ML threat detection profiles
config/cato/casb-genai.yaml	Generative AI CASB policies

## 57. AGI Brain & Identity Data Fabric: Agentic Orchestration

The final layer of the RADIANT stack is the AGI Brain, powered by Radiant Logic and the Identity Data Fabric. This layer transforms raw compute and connectivity into intelligent action, moving beyond static "automation" to dynamic "agentic" behavior.

### 57.1 The Identity Data Fabric

Radiant Logic pioneered the **Identity Data Fabric**, a unified layer that abstracts the complexity of identity data across the enterprise. In an AGI environment, "identity" is not just for humans—it includes agents, APIs, microservices, and even the Genesis reactors themselves.

#### Identity Types in RADIANT

Identity Type	Description	Authentication Method
<b>Human Users</b>	End users of Think Tank	Cognito + MFA
<b>Administrators</b>	Platform operators	Cognito + Hardware Key
<b>AI Agents</b>	Autonomous software agents	Service Account + JWT
<b>API Clients</b>	External integrations	API Key + IP Allowlist
<b>Microservices</b>	Internal services	mTLS + Service Mesh
<b>Genesis Reactors</b>	Physical infrastructure	mTLS + Hardware HSM
<b>Sentinel Agents</b>	Background monitoring agents	Ephemeral tokens

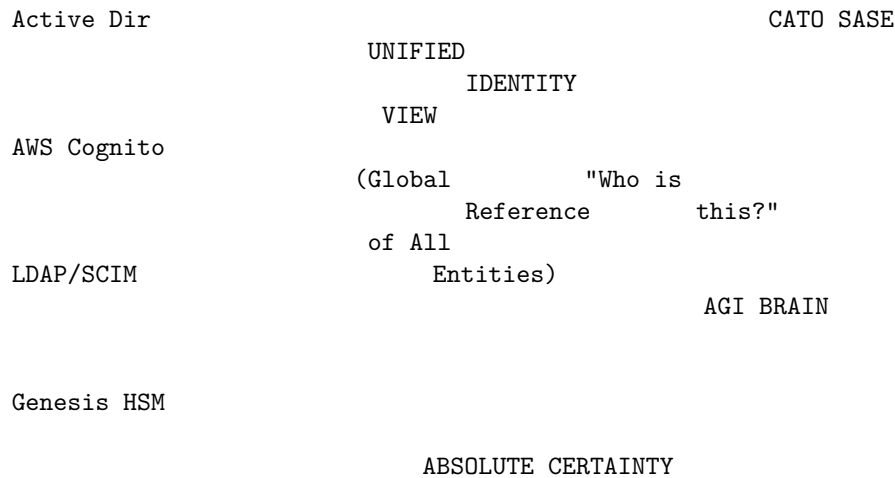
**The RadiantOne Platform** The RadiantOne Platform serves as the **central nervous system**, maintaining a global reference of all entities. This is crucial for the Zero Trust model enforced by Cato:

#### IDENTITY DATA FABRIC ARCHITECTURE

IDENTITY SOURCES

RADIANTONE

CONSUMERS



The Cato network asks **"Who are you?"** and the Radiant Brain answers with **absolute certainty**, backed by real-time data from the Identity Fabric.

## 57.2 Agentic AI and the Radiant Ghost

The transition from traditional scripting to **Agentic AI** is a key architectural shift. Unlike automation, which follows a rigid set of if-then rules, AI agents continuously evaluate their environment, learn from outcomes, and adapt their behavior.

**The Radiant Ghost Metaphor** To visualize agentic behavior for users, the interface employs the metaphor of the **"Radiant Ghost"**. This is not a spooky or malicious entity, but a **benevolent, semi-autonomous agent** that "haunts" the network to protect it.

Ghost State	Visual Indicator	Meaning
<b>Dormant</b>	Faint outline	Agent is monitoring but not acting
<b>Active</b>	Pulsing glow	Agent is actively processing a task
<b>Hunting</b>	Searching animation	Agent is investigating a potential threat
<b>Remediating</b>	Repair animation	Agent is autonomously fixing an issue
<b>Alerting</b>	Red pulse	Agent requires human attention

**Ghost Vector Integration** The "Ghost" vector serves as the UI avatar for background processes. When an administrator sees the glowing "Radiant Brain" or "Ghost" icon, they know that the agent is actively:

- Hunting threats across the network
- Optimizing performance based on telemetry
- Scrubbing orphan identities from the directory
- Consolidating memories in the Grimoire
- Monitoring Genesis reactor telemetry

## 57.3 Technical Standards: MCP, SSF, and CAEP

For proper configuration, administrators must understand the specific protocols used by AI agents.

**Model Context Protocol (MCP)** The system leverages the **Model Context Protocol (MCP)** to standardize how AI agents interface with the underlying data models. This ensures that an agent analyzing a security log understands the **context** of that log within the broader Identity Fabric.

```
// MCP configuration for identity-aware agents
export const MCP_IDENTITY_CONFIG = {
  // Context providers
  contextProviders: [
    {
      name: 'identity_fabric',
      endpoint: 'radiantone://identity/v1',
      capabilities: ['user_lookup', 'group_membership', 'entitlements'],
    },
    {
      name: 'genesis_telemetry',
      endpoint: 'genesis://telemetry/v1',
      capabilities: ['power_status', 'safety_status', 'maintenance_schedule'],
    },
    {
      name: 'cato_security',
      endpoint: 'cato://security/v1',
      capabilities: ['threat_status', 'access_policies', 'blocked_entities'],
    },
  ],

  // Tool definitions for agents
  tools: [
    {
      name: 'lookup_identity',
      description: 'Look up identity information from the fabric',
      parameters: {
        identifier: 'string',
        identifier_type: 'user_id | email | service_account | device_id',
      },
    },
    {
      name: 'remediate_orphan',
      description: 'Remove orphan account from directory',
      parameters: {
        account_id: 'string',
        reason: 'string',
      },
      requires_approval: false, // Autonomous remediation enabled
    },
    {
      name: 'escalate_to_human',
      description: 'Escalate decision to human operator',
      parameters: {
        severity: 'low | medium | high | critical',
        context: 'object',
        recommended_action: 'string',
      },
    },
  ],
}
```

```
};
```

**Shared Signals Framework (SSF) and CAEP** See Section 56.4 for detailed SSF/CAEP configuration. The AGI Brain consumes CAEP signals from:

- Genesis Reactor (physical events)
- Cato SASE (network events)
- Radiant Logic (identity events)

## 57.4 fastWorkflow: The Logic Engine

The reliability of AI agents is secured by **fastWorkflow**, an open-source framework designed for building complex, large-scale Python applications. AI agents can sometimes be unpredictable ("hallucinations"). fastWorkflow provides:

Capability	Description
<b>Robust Validation Pipeline</b>	Every agent output is validated before execution
<b>Deterministic Business Logic</b>	Core decision paths are pre-programmed and immutable
<b>Proper Error Handling</b>	Failures are caught, logged, and escalated appropriately
<b>Audit Logging</b>	Every action is recorded for compliance
<b>Reliable Tool Execution</b>	When an agent decides to "isolate a host," it executes correctly

This "Reliable Tool Execution" is **critical** when the agent is controlling nuclear-powered infrastructure.

### fastWorkflow Configuration

```
# fastworkflow_config.py
from fastworkflow import Workflow, Task, ValidationPipeline

# Define validation pipeline for agent actions
validation_pipeline = ValidationPipeline([
    # Pre-execution validation
    ("syntax_check", lambda action: action.is_syntactically_valid()),
    ("permission_check", lambda action: action.has_required_permissions()),
    ("safety_check", lambda action: action.passes_genesis_interlock()),

    # Post-execution validation
    ("result_validation", lambda result: result.matches_expected_schema()),
    ("side_effect_check", lambda result: result.side_effects_are_acceptable()),
])

# Define deterministic safety constraints (cannot be overridden by AI)
IMMUTABLE_CONSTRAINTS = {
    "genesis_interlock": True,
    "require_human_approval_for_destructive_actions": True,
    "max_autonomous_remediations_per_hour": 100,
    "blocked_actions_during_genesis_alert": [
        "delete_identity",
        "revoke_all_access",
        "shutdown_service",
    ],
}

# Agent workflow definition
```

```

class IdentityRemediationWorkflow(Workflow):
    """Autonomous identity remediation with fastWorkflow safety"""

    @Task(validation=validation_pipeline)
    def identify_orphan_accounts(self, criteria: dict) -> list:
        """Identify orphan accounts in the Identity Fabric"""
        # Deterministic query logic
        pass

    @Task(validation=validation_pipeline, requires_approval=False)
    def remediate_orphan(self, account_id: str, reason: str) -> dict:
        """Remove orphan account (autonomous, no approval needed)"""
        # Pre-check: Genesis Interlock
        if not self.genesis_interlock_allows_action():
            raise GenesisInterlockException("Action blocked during Genesis alert")

        # Execute remediation
        pass

    @Task(validation=validation_pipeline, requires_approval=True)
    def bulk_remediation(self, account_ids: list) -> dict:
        """Bulk remediation requires human approval"""
        pass

```

## 57.5 Autonomous Identity Remediation

The RadiantOne platform can be configured to allow **Autonomous Remediation**. This grants AI agents the permission to fix data quality issues (e.g., orphan accounts) without human intervention.

**Enabling Autonomous Remediation** Navigate to **Identity → Remediation → Autonomous Settings**:

```

# autonomous-remediation.config.yaml
autonomous_remediation:
    enabled: true

# Actions that can be performed without human approval
autonomous_actions:
    - action: remove_orphan_account
      conditions:
        - account_inactive_days: 90
        - no_associated_entitlements: true
        - no_recent_authentication: true
      max_per_hour: 50

    - action: disable_stale_service_account
      conditions:
        - last_used_days: 180
        - not_in_critical_systems: true
      max_per_hour: 20

    - action: fix_group_membership_inconsistency
      conditions:
        - source_of_truth_mismatch: true
      max_per_hour: 100

```

```

# Actions that ALWAYS require human approval
always_require_approval:
  - delete_human_account
  - revoke_admin_privileges
  - modify_genesis_access
  - bulk_operations_over_100

# Genesis Interlock integration
genesis_interlock:
  pause_during_alert_levels: [GEN-300, GEN-400, GEN-500]
  resume_automatically: true

```

## 57.6 The Memory Safety Imperative

The Think Tank identifies **memory safety vulnerabilities** as a specific class of defect responsible for **70% of all software vulnerabilities**. These bugs allow attackers to alter data and command systems—a catastrophic risk for a nuclear-powered AGI.

**AI-Driven Code Refactoring** The agents within the RADIANT ecosystem are not just managing identity; they can actively refactor the codebase of the infrastructure itself to eliminate memory safety errors. This fulfills the Genesis mandate for flawlessness.

Unsafe Language	Safe Alternative	AI Capability
C	Rust	Automated translation of security-critical components
C++	Rust	Automated translation with human review
Assembly	Safe abstractions	Identification and encapsulation

## Memory Safety Scanning Configuration

```

# memory-safety.config.yaml
memory_safety:
  enabled: true

# Automated scanning
scanning:
  schedule: daily
  targets:
    - path: /src/genesis-interface/**
      priority: critical
    - path: /src/network/**
      priority: high
    - path: /src/ai-agents/**
      priority: medium

# AI-assisted refactoring
ai_refactoring:
  enabled: true
  auto_submit_pr: true
  require_human_review: true
  target_language: rust

# Vulnerability classes to detect

```



```
vulnerability_classes:
- buffer_overflow
- use_after_free
- double_free
- null_pointer_dereference
- integer_overflow
- format_string_vulnerability
```

57.7 Database Tables

Table	Purpose
identity_fabric_sync_log	Sync history with identity sources
autonomous_remediation_log	Audit trail of autonomous actions
agent_activity_log	All AI agent activities
mcp_context_queries	MCP context provider queries
memory_safety_scans	Code scanning results
memory_safety_remediations	AI-assisted code fixes

57.8 Implementation Files

File	Purpose
lambda/shared/services/identity-fabric.service.ts	Identity Data Fabric integration
lambda/shared/services/autonomous-remediation.service.ts	Autonomous remediation engine
lambda/shared/services/mcp-identity-provider.service.ts	MCP context for identity
lambda/shared/services/memory-safety-scanner.service.ts	Code safety scanning
python/cato/agents/identity_remediation.py	fastWorkflow identity agent
migrations/163_agi_brain_identity.sql	Database schema
config/identity/autonomous-remediation.yaml	Remediation configuration

58. Deployment Safety & Environment Management

58.1 Overview

RADIANT uses a **three-environment architecture** (Dev, Staging, Prod) with strict safety rules to prevent accidental infrastructure damage. This section documents the **hard rules** that must never be bypassed.

58.2 Environment Architecture

Environment	AWS Profile	Purpose	CDK Watch	Approval Required
Dev	radiant-dev	Development, testing	Allowed	No
Staging	radiant-staging	Pre-production testing	<b>FORBIDDEN</b>	Yes
Prod	radiant-prod	Production	<b>FORBIDDEN</b>	Yes

58.3 Critical Safety Rule: cdk watch is DEV-ONLY

THIS RULE MUST NEVER BE IGNORED

cdk watch --hotswap is **ONLY** allowed in the DEV environment. It is **absolutely forbidden** for staging and production.

**Why This Rule Exists** `cdk watch --hotswap` bypasses CloudFormation safety mechanisms:

Risk	Description	Impact
No Rollback	Hotswap changes cannot be rolled back automatically	Manual recovery required
State Drift	CloudFormation state doesn't match actual resources	Future deployments may fail
Inconsistent Infrastructure	Partial deployments can leave broken state	Service outages
No Change Sets	No preview of what will change	Unexpected deletions

**Enforcement Points** The rule is enforced at **four independent levels**:

1. **Swift Deployer App** (apps/swift-deployer/Sources/RadiantDeployer/Services/CDKService.swift)

```
guard isHotswapAllowed(environment: environment) else {  
    throw CDKError.hotswapBlockedForEnvironment(environment)  
}
```

- `deployWithHotswap()` and `startWatch()` throw errors for non-dev
- Standard `deploy()` uses `--require-approval` broadening for staging/prod
- **This is the primary deployment method and enforces safety automatically**

2. **CDK Entry Point** (packages/infrastructure/bin/radiant.ts)

```
if (isCdkWatch && detectedEnv !== 'dev') {  
    console.error(' BLOCKED: cdk watch is FORBIDDEN...');  
    process.exit(1); // Hard exit - no bypass  
}
```

3. **Environment Configuration** (packages/infrastructure/lib/config/environments.ts)

- `enableCdkWatch: false` hardcoded for staging/prod
- `assertCdkWatchAllowed(env)` throws error for non-dev

4. **Safety Script** (scripts/cdk-safety-check.sh)

- Pre-deploy check that blocks watch on non-dev
- Requires typing environment name to confirm staging/prod deploys

## 58.4 Safe Deployment Methods

**For Development** (`cdk watch` allowed)

```
# Configure credentials  
source ./scripts/setup_credentials.sh  
  
# Start Direct Dev Mode  
cd packages/infrastructure  
npx cdk watch --hotswap --profile radiant-dev
```

**For Staging** (approval required)

```
# Option 1: Swift Deployer (recommended)  
# Open Swift Deployer app → Select Staging → Deploy  
  
# Option 2: CLI with approval gates  
AWS_PROFILE=radiant-staging npx cdk deploy --all \  
--require-approval broadening
```

## For Production (approval required)

```
# Option 1: Swift Deployer (STRONGLY recommended)
# Open Swift Deployer app → Select Production → Deploy

# Option 2: CLI with approval gates (use with caution)
AWS_PROFILE=radiant-prod npx cdk deploy --all \
  --require-approval broadening
```

## 58.5 Credential Setup

Before deploying to any environment, configure AWS credentials:

1. **Create IAM Users** in AWS Console:

- radiant-dev-deployer
- radiant-staging-deployer
- radiant-prod-deployer

2. **Attach Permissions:** AdministratorAccess policy (or scoped CDK policy)

3. **Generate Access Keys** for each user

4. **Configure Profiles:**

```
# Edit the credentials script first
nano scripts/setup_credentials.sh

# Then run it
source ./scripts/setup_credentials.sh
```

5. **Bootstrap CDK** (one-time per account/region):

```
chmod +x ./scripts/bootstrap_cdk.sh
./scripts/bootstrap_cdk.sh
```

## 58.6 Environment Detection

The CDK entry point detects the target environment via multiple signals:

Signal	Priority	Example
RADIANT_ENV env var	1 (highest)	RADIANT_ENV=staging
CDK_CONTEXT_environment	2	-c environment=staging
AWS_PROFILE pattern	3	radiant-staging in profile name
Default	4 (lowest)	Falls back to dev

## 58.7 Troubleshooting

**”BLOCKED: cdk watch is FORBIDDEN” Error** This error means you attempted to run `cdk watch` against a non-dev environment. This is intentional and cannot be bypassed.

**Solution:** Use `cdk deploy` with approval gates instead:

```
AWS_PROFILE=radiant-{env} npx cdk deploy --all --require-approval broadening
```

**Wrong Environment Detected** If the wrong environment is being detected:

```
# Explicitly set the environment
export RADIANT_ENV=dev
export AWS_PROFILE=radiant-dev

# Verify
echo $RADIANT_ENV $AWS_PROFILE
```

## 58.8 Implementation Files

File	Purpose
apps/swift-deployer/.../CDKService.swift	Swift Deployer CDK service with safety enforcement
packages/infrastructure/bin/radiant.ts	CDK entry point with safety check
packages/infrastructure/lib/config/environments.ts	Environment configurations
scripts/setup_credentials.sh	AWS credential setup
scripts/bootstrap_cdk.sh	CDK bootstrap helper
scripts/cdk-safety-check.sh	Pre-deploy safety validation
packages/infrastructure/ARCHITECTURE.md	Stack vs Lambda architecture guide

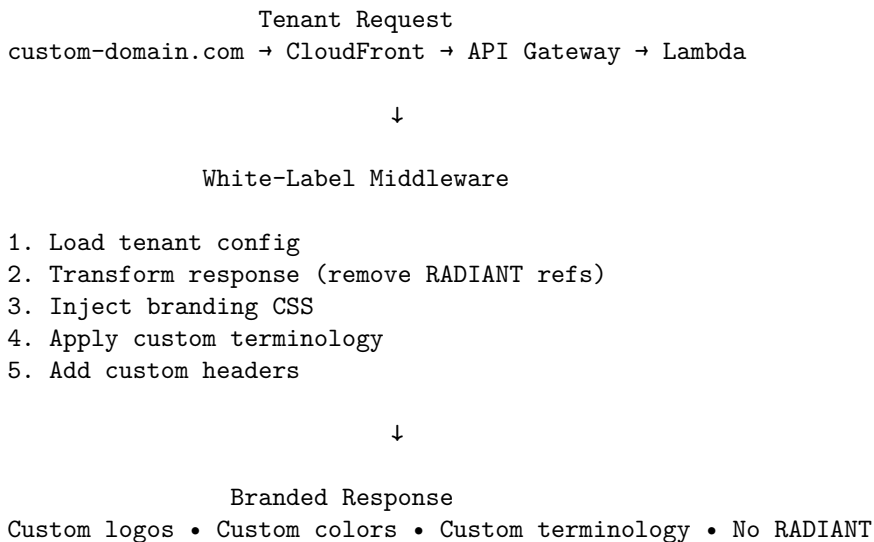
## 59. White-Label Invisibility (Moat #25)

**Moat Evaluation:** Score 18/30 - Tier 4 Business Model Moat. End users never know RADIANT exists. Infrastructure stickiness creates platform layer dependency.

### 59.1 Overview

White-Label Invisibility enables tenants to completely rebrand the platform so their end users never see RADIANT branding. This creates strong infrastructure stickiness—switching providers means rebuilding their entire branded experience.

### 59.2 Architecture



### 59.3 Branding Configuration

Setting	Type	Description
companyName	string	Company name displayed throughout
productName	string	Product name (replaces "Think Tank")
tagline	string	Optional tagline
logo.primary	URL	Main logo URL
logo.light	URL	Light theme logo
logo.dark	URL	Dark theme logo
logo.icon	URL	Favicon/icon
colors.primary	hex	Primary brand color
colors.secondary	hex	Secondary color
colors.accent	hex	Accent color
fonts.primary	string	Primary font family
fonts.secondary	string	Secondary font family
fonts.mono	string	Monospace font

### 59.4 Feature Visibility

Control what end users can see:

Setting	Default	Description
hideRadiantBranding	true	Remove all RADIANT references
hidePoweredBy	true	Hide "Powered by RADIANT" footer
hideModelNames	false	Anonymize model names to "AI Model"
hideModelProviders	false	Remove provider info (Anthropic, OpenAI)
hideCostMetrics	false	Hide cost information from UI
hideUsageMetrics	false	Hide usage statistics
disabledFeatures	[]	Features to disable for this tenant

### 59.5 Custom Terminology

Replace standard terms with branded alternatives:

Default Term	Example Custom
AI Assistant	"Alex" (custom assistant name)
Conversation	"Session"
Artifact	"Creation"
Workspace	"Studio"
Team	"Organization"

### 59.6 Custom Domains

#### Adding a Custom Domain

1. Add domain in admin dashboard
2. Get DNS verification record
3. Add TXT record to DNS
4. Wait for verification (auto-checks every hour)
5. SSL certificate auto-provisioned via ACM
6. CloudFront distribution updated

## Domain Types

Type	Purpose
<code>primary</code>	Main application domain
<code>alias</code>	Additional domains that redirect
<code>api</code>	API-specific subdomain

## 59.7 Email Templates

Customize all outbound emails:

Template	Purpose
<code>welcome</code>	New user welcome email
<code>password_reset</code>	Password reset link
<code>invitation</code>	Team invitation
<code>notification</code>	General notifications
<code>billing</code>	Billing-related emails

Each template supports:

- Custom subject line
- HTML template with variables
- Plain text fallback
- Custom from name/email

## 59.8 Legal Configuration

Setting	Description
<code>companyLegalName</code>	Legal entity name
<code>termsOfServiceUrl</code>	Link to ToS
<code>privacyPolicyUrl</code>	Link to privacy policy
<code>cookiePolicyUrl</code>	Link to cookie policy
<code>supportEmail</code>	Support contact email
<code>copyrightNotice</code>	Footer copyright text
<code>customFooterHtml</code>	Optional custom footer HTML

## 59.9 API Customization

Setting	Description
<code>customBaseUrl</code>	Custom API base URL
<code>hideVersionHeader</code>	Remove X-Radiant-Version header
<code>customHeaders</code>	Additional headers to inject
<code>corsOrigins</code>	Allowed CORS origins
<code>rateLimitOverrides</code>	Per-endpoint rate limit overrides

## 59.10 Response Transformation

The white-label middleware automatically transforms responses:

```
// Before transformation
{
  "model": "claude-3.5-sonnet",
  "provider": "anthropic",
  "message": "Welcome to RADIANT Think Tank!"
}

// After transformation (with hideModelNames + hideModelProviders)
{
  "model": "AI Model",
  "message": "Welcome to [ProductName]!"
}
```

### 59.11 CSS Injection

Custom CSS is automatically injected for brand consistency:

```
:root {
  --color-primary: #3B82F6;      /* From config */
  --color-secondary: #6366F1;
  --font-primary: 'Inter', sans-serif;
  /* ... all brand tokens */
}
```

### 59.12 API Endpoints

Endpoint	Method	Description
/api/admin/white-label/config	GET	Get configuration
/api/admin/white-label/config	POST	Create configuration
/api/admin/white-label/config	PUT	Update configuration
/api/admin/white-label/config	DELETE	Delete configuration
/api/admin/white-label/validate	POST	Validate configuration
/api/admin/white-label/domains	POST	Add domain
/api/admin/white-label/domains/:id	DELETE	Remove domain
/api/admin/white-label/domains/:domain/verify	POST	Initiate verification
/api/admin/white-label/domains/:domain/verify	GET	Check verification
/api/admin/white-label/branding	PUT	Update branding
/api/admin/white-label/features	PUT	Update feature visibility
/api/admin/white-label/legal	PUT	Update legal config
/api/admin/white-label/emails	PUT	Update email config
/api/admin/white-label/preview	GET	Generate branding preview
/api/admin/white-label/export	GET	Export configuration
/api/admin/white-label/import	POST	Import configuration
/api/admin/white-label/metrics	GET	Get usage metrics

### 59.13 Database Tables

Table	Purpose
white_label_config	Per-tenant configuration
white_label_domains	Custom domains
domain_verifications	DNS verification records
branding_assets	Uploaded logos, fonts, images

Table	Purpose
white_label_email_templates	Custom email templates
custom_terminology	Term mappings
white_label_metrics	Usage metrics

## 59.14 Implementation Files

File	Purpose
packages/shared/src/types/white-label.types.ts	Type definitions
lambda/shared/services/white-label.service.ts	Core service
lambda/admin/white-label.ts	API handler
migrations/172_white_label.sql	Database schema

## 59.15 Usage Example

```
import { whiteLabelService } from './services/white-label.service';

// Create white-label configuration
const config = await whiteLabelService.createConfig(tenantId, {
  enabled: true,
  branding: {
    companyName: 'Acme Corp',
    productName: 'Acme AI',
    logo: { primary: 'https://...', light: '...', dark: '...', icon: '...' },
    colors: { primary: '#FF5733', secondary: '#3366FF', ... },
  },
  features: {
    hideRadiantBranding: true,
    hidePoweredBy: true,
    hideModelNames: true,
  },
  legal: {
    companyLegalName: 'Acme Corporation Inc.',
    supportEmail: 'support@acme.com',
    copyrightNotice: '© 2026 Acme Corp',
  },
});

// Add custom domain
await whiteLabelService.addDomain(tenantId, 'ai.acme.com', 'primary');

// Transform API responses
const transformedResponse = whiteLabelService.transformResponse(tenantId, originalResponse);
```

## 60. User Violation Enforcement System

The User Violation Enforcement System provides comprehensive tracking, escalation, and enforcement of regulatory and policy violations. This is critical for HIPAA, GDPR, and SOC2 compliance.



## 60.1 Overview

The system enables administrators to:

- Report and track policy/regulatory violations per user
- Configure automatic escalation policies
- Take enforcement actions (warnings through account termination)
- Process user appeals with audit trail
- Monitor high-risk users with risk scoring

## 60.2 Violation Categories

Category	Description	Examples
hipaa	HIPAA violations	PHI exposure, unauthorized access, sharing
gdpr	GDPR violations	Consent violations, retention issues, cross-border transfers
soc2	SOC2 violations	Security control failures
terms_of_service	ToS violations	Terms breach
acceptable_use	AUP violations	Misuse of platform
content_policy	Content violations	Harmful, illegal content
security	Security violations	Credential sharing, injection attempts
billing	Billing violations	Payment fraud, chargeback abuse
abuse	Platform abuse	Rate limit abuse, spam

## 60.3 Severity Levels

Severity	Description	Risk Score	Impact
warning	Minor issue, first offense	+5	
minor	Low impact violation	+10	
major	Significant violation	+20	
critical	Severe violation requiring immediate action	+40	

## 60.4 Enforcement Actions

Action	Description
warning_issued	Formal warning recorded
feature_restricted	Specific features disabled
rate_limited	API/usage rate limits applied
temporarily_suspended	Account suspended for specified duration
permanently_suspended	Account permanently suspended
account_terminated	Account deleted
reported_to_authorities	Reported to regulatory authorities

## 60.5 Configuration

```
interface UserViolationConfig {
    enabled: boolean;           // Enable violation tracking
    autoDetectionEnabled: boolean; // Auto-detect from system events
    autoEnforcementEnabled: boolean; // Auto-apply escalation actions

    // Notifications
    notifyUserOnViolation: boolean; // Notify user when violation reported
}
```

```

notifyUserOnAction: boolean;           // Notify user when action taken
notifyAdminOnCritical: boolean;        // Alert admins on critical violations
adminNotificationEmails: string[];     // Admin email addresses

// Retention
retentionDays: number;                 // Default: 2555 (~7 years)

// Appeals
allowAppeals: boolean;                // Allow users to appeal
appealWindowDays: number;              // Days to submit appeal (default: 30)
maxAppealsPerViolation: number;        // Max appeals per violation (default: 2)

// Compliance
requireEvidenceRedaction: boolean;     // Always redact evidence content
auditAllActions: boolean;              // Log all actions to audit trail
}

```

## 60.6 Escalation Policies

Configure automatic enforcement based on violation patterns:

```

interface EscalationRule {
  triggerType: 'count' | 'severity' | 'category';

  // Count-based: trigger after N violations
  violationCount?: number;
  withinDays?: number;

  // Severity-based: trigger on severity threshold
  severityThreshold?: ViolationSeverity;

  // Category-based: trigger for specific categories
  categories?: ViolationCategory[];

  // Action to take
  action: EnforcementAction;
  actionDurationDays?: number;           // For temporary actions
  requiresManualReview: boolean;         // Require admin review
  allowAppeal: boolean;                  // Allow user appeal
}

```

### Example Policy:

- 3 warnings in 90 days → Feature restriction
- 2 major violations → Temporary suspension (7 days)
- 1 critical violation → Immediate suspension (requires review)
- Any HIPAA violation → Immediate suspension + admin alert

## 60.7 API Endpoints

Base URL: /api/admin/violations

Endpoint	Method	Description
/dashboard	GET	Get dashboard data with metrics
/config	GET	Get configuration

Endpoint	Method	Description
/config	PUT	Update configuration
/violations	GET	Search violations
/violations	POST	Report new violation
/violations/:id	GET	Get violation details
/violations/:id	PUT	Update violation
/violations/:id/action	POST	Take enforcement action
/users/:userId/violations	GET	Get user's violations
/users/:userId/summary	GET	Get user violation summary
/users/:userId/suspend	POST	Suspend user
/users/:userId/reinstate	POST	Reinstate user
/appeals	GET	Get pending appeals
/appeals/:id	GET	Get appeal details
/appeals/:id/review	POST	Review appeal
/metrics	GET	Get violation metrics
/policies	GET	Get escalation policies
/policies	POST	Create escalation policy

## 60.8 Admin Dashboard

Access at: **Compliance** → **Violations** (/compliance/violations)

### Features:

- **Dashboard:** Metrics cards (total, new, resolved, pending appeals, avg resolution time)
- **Violations Tab:** Search/filter violations, report new, take actions
- **Appeals Tab:** Review and decide on pending appeals
- **High Risk Users Tab:** View users with elevated risk scores
- **Settings Tab:** Configure system settings

## 60.9 User Appeal Process

1. User submits appeal within appeal window (default: 30 days)
2. Violation status changes to **appealed**
3. Admin reviews appeal
4. Decisions:
  - **Upheld:** Original action stands
  - **Overtured:** Violation dismissed, action reversed
  - **Reduced:** Lesser action applied
5. User notified of decision
6. All actions logged to audit trail

## 60.10 Risk Scoring

Users are assigned a risk level based on active violations:

Risk Level	Criteria
low	0-1 active violations
moderate	2+ active violations
elevated	5+ active violations
high	2+ major violations active
critical	Any critical violation active

Risk score (0-100) is calculated as sum of severity impacts for active violations.

## 60.11 Database Tables

Table	Purpose
<code>user_violations</code>	Violation records with enforcement
<code>violation_evidence</code>	Evidence attachments (redacted)
<code>violation_appeals</code>	Appeal records with review
<code>violation_escalation_policies</code>	Escalation policy definitions
<code>violation_escalation_rules</code>	Policy rules/thresholds
<code>user_violation_config</code>	Per-tenant configuration
<code>user_violation_summary</code>	Aggregated user summaries
<code>violation_audit_log</code>	Immutable audit trail

## 60.12 Implementation Files

File	Purpose
<code>packages/shared/src/types/user-violations.types.ts</code>	Type definitions
<code>lambda/shared/services/user-violation.service.ts</code>	Core service
<code>lambda/admin/user-violations.ts</code>	API handler
<code>apps/admin-dashboard/app/(dashboard)/compliance/violations/page.tsx</code>	Admin UI
<code>migrations/173_user_violations.sql</code>	Database schema

## 60.13 Compliance Considerations

- **HIPAA:** 7-year retention, audit trail, PHI redaction in evidence
- **GDPR:** User notification, appeal rights, data minimization
- **SOC2:** Audit logging, access controls, incident response

# 61. Multi-Protocol Gateway Architecture

## 61.1 Overview

The RADIANT Multi-Protocol Gateway v3.0 is a custom Go-based connectivity layer supporting MCP, A2A, OpenAI, Anthropic, and Google interfaces at **1M+ concurrent connection scale**.

### Key Components:

Component	Technology	Purpose
<b>Go Gateway</b>	Custom Go service	WebSocket/SSE termination, NATS bridging
<b>NATS JetStream</b>	Message bus	At-least-once delivery, session persistence
<b>Cedar Authorization</b>	ABAC policies	Resource-level access control
<b>Resume Tokens</b>	HMAC-signed tokens	Session rehydration on reconnect

## 61.2 Architecture Diagram

Internet → NLB (Layer 4) → Go Gateway Fleet → NATS JetStream → Lambda Workers  
→ ~80K connections per c6g.xlarge instance

### 61.3 Protocol Support

Protocol	Subject Pattern	Use Case
<b>MCP</b>	<code>in.mcp.{tenant}.{agent}</code>	Model Context Protocol
<b>A2A</b>	<code>in.a2a.{tenant}.{agent}</code>	Agent-to-Agent communication
<b>FC</b>	<code>in.fc.{tenant}.{agent}</code>	Function calling

### 61.4 NATS Stream Configuration

Stream	Subjects	Retention	Purpose
INBOX	<code>in.&gt;</code>	WorkQueue	Incoming messages
OUTBOX	<code>out.&gt;</code>	1h TTL	Responses to sessions
HISTORY	<code>history.&gt;</code>	10K msgs/subject	Session replay

### 61.5 Cedar Authorization

Resource-level ABAC with policies for:

- Cross-tenant access denial
- Tool namespace restrictions
- Sensitive resource protection
- Admin bypass rules

### 61.6 Resume Token Strategy

Sessions survive connection drops via HMAC-signed resume tokens containing:

- Session ID, Tenant ID, Principal ID
- NATS inbox/outbox subjects
- Expiry timestamp

### 61.7 Capacity Planning

Component	Instance	Capacity	Count for 1M
Go Gateway	c6g.xlarge	80K conn	13 instances
NATS Cluster	r6g.xlarge	N/A	3 nodes
Lambda (MCP)	1024MB	N/A	1000 concurrent

**Estimated Cost:** \$8,000-15,000/month for 1M connections

### 61.8 Admin Dashboard Controls

The Gateway admin interface provides comprehensive monitoring and configuration:

#### Dashboard Overview (/gateway)

- Real-time connection count and peak connections
- Messages per minute throughput
- Average latency with warning thresholds
- Error rate monitoring
- Active instance count

#### Statistics & Reporting

- 5-minute bucketed statistics stored persistently
- Hourly and daily aggregated views
- Protocol distribution (MCP, A2A, OpenAI, Anthropic, Google)
- 24-hour trend visualization

### Configuration Controls

- Connection limits (per tenant, per user, per agent)
- Rate limits (messages/second, bytes/second)
- Timeout settings (connect, idle, read, write)
- Protocol enable/disable toggles
- mTLS enforcement for A2A
- Resume token TTL configuration

### Maintenance Mode

- Enable/disable with custom message
- IP allowlist for maintenance bypass
- Graceful connection draining

### Alert Management

- Severity levels: info, warning, critical
- Alert types: connection\_limit, rate\_limit, error\_spike, latency\_spike, instance\_unhealthy
- Acknowledge and resolve workflows
- Resolution notes tracking

### Instance Management

- View all gateway instances with status
- Drain instances gracefully
- Monitor per-instance metrics

**API Endpoints** (/api/admin/gateway) | Endpoint | Method | Purpose | |-----|-----|-----| |

/dashboard	GET	Overview metrics and alerts		/statistics	GET	Time-series statistics	
/configuration	GET/PUT	View/update configuration		/maintenance	POST	Set maintenance mode	
/alerts	GET	List alerts		/alerts/:id/acknowledge	POST	Acknowledge alert	
/alerts/:id/resolve	POST	Resolve alert		/sessions	GET	List active sessions	
/sessions/:id/terminate	POST	Terminate session		/instances	GET	List gateway instances	
/instances/:id/drain	POST	Drain instance					

## 61.9 Database Schema

### Tables Created:

- gateway\_instances - Instance registry with heartbeat tracking
- gateway\_statistics - Time-series metrics (5-minute buckets)
- gateway\_configuration - Per-tenant and global settings
- gateway\_alerts - Alert and incident tracking
- gateway\_sessions - Active connection tracking
- gateway\_audit\_log - Admin action audit trail

### Report Types Added:

- gateway-statistics - Connection and message statistics
- gateway-alerts - Alert summary report

## 61.10 Implementation Files

File	Purpose
apps/gateway/	Go Gateway service (16 files)
apps/gateway/internal/server/tls.go	TLS/mTLS configuration
services/egress-proxy/	HTTP/2 connection pool
infrastructure/cedar/schema.cedarschema	Cedar entity/action schema
infrastructure/cedar/policies/	Authorization policies
packages/infrastructure/lambda/shared/services/cedar/	Cedar TypeScript service
packages/infrastructure/lambda/gateway/mcp-worker.ts	MCP NATS consumer Lambda
packages/infrastructure/lambda/admin/gateway.ts	Gateway admin API Lambda
packages/infrastructure/__tests__/gateway/	Gateway integration tests
packages/infrastructure/migrations/V2026_01_20_001__gateway_statistics.sql	Statistics schema
apps/admin-dashboard/app/(dashboard)/gateway/page.tsx	Admin dashboard UI
apps/thinktank-admin/app/(dashboard)/gateway/page.tsx	Think Tank status view
infrastructure/docker/gateway/	Local dev stack
infrastructure/load-tests/	k6 load testing scripts
packages/infrastructure/lib/stacks/gateway-stack.ts	CDK deployment

61.11 Documentation Reference

For complete architecture details, see:

- **MULTI-PROTOCOL-GATEWAY-ARCHITECTURE.md** - Full technical specification

Section 62: Code Quality & Test Coverage Visibility

62.1 Overview

The Code Quality dashboard provides real-time visibility into:

- **Test Coverage:** Line, function, and branch coverage by component
- **Technical Debt:** Tracking items aligned with `TECHNICAL_DEBT.md`
- **JSON Safety Migration:** Progress migrating `JSON.parse` to safe utilities
- **Code Quality Alerts:** Automated alerts for coverage drops and regressions

62.2 Admin Dashboard

**Location:** `/code-quality`

The Code Quality page provides:

- **Summary Cards:** Overall coverage %, open debt items, JSON safety progress, active alerts
- **Coverage by Component:** Breakdown for lambda, admin-dashboard, swift-deployer
- **Technical Debt Table:** Prioritized list with status, estimated hours
- **JSON Safety Progress:** Migration status per component with progress bars
- **Alerts Tab:** Active code quality alerts with acknowledge/resolve actions

62.3 API Endpoints

**Base Path:** `/api/admin/code-quality`

Endpoint	Method	Purpose
<code>/dashboard</code>	GET	Overview metrics and summary
<code>/coverage</code>	GET	Latest coverage by component
<code>/coverage/history</code>	GET	Coverage trend over time

Endpoint	Method	Purpose
/debt	GET	List technical debt items
/debt/:id	PUT	Update debt item status
/json-safety	GET	JSON migration progress
/json-safety/locations	GET	List JSON.parse locations
/alerts	GET	List code quality alerts
/alerts/:id/acknowledge	POST	Acknowledge an alert
/alerts/:id/resolve	POST	Resolve an alert
/files-needing-tests	GET	Files that need test coverage

## 62.4 Database Schema

Tables Created (Migration V2026\_01\_20\_002\_\_code\_quality\_metrics.sql):

Table	Purpose
code_quality_snapshots	Periodic snapshots of coverage/quality metrics
test_file_registry	Registry of source files and test status
json_parse_locations	Tracking JSON.parse for migration
technical_debt_items	Debt items aligned with TECHNICAL_DEBT.md
code_quality_alerts	Alerts for quality regressions

Views:

- v\_latest\_test\_coverage - Latest coverage per component
- v\_json\_safety\_progress - JSON migration progress
- v\_technical\_debt\_summary - Debt summary by priority

## 62.5 Report Integration

The reporting system includes a `code_quality` report type:

```
const report = await reportGeneratorService.generateReport({
  id: 'code-quality-report',
  tenant_id: tenantId,
  name: 'Code Quality Report',
  report_type: 'code_quality',
  format: 'pdf',
  parameters: {},
  recipients: ['admin@example.com'],
});
```

Report Sections:

1. Test Coverage by Component
2. Technical Debt Items
3. JSON Safety Migration Progress

## 62.6 Safe JSON Utilities

Located in `lambda/shared/utils/safe-json.ts`:

Function	Purpose
<code>safeJsonParse(json, fallback?)</code>	Parse with optional fallback



Function	Purpose
<code>parseJsonWithSchema(json, schema)</code>	Parse with Zod validation
<code>parseJsonWithSchemaOrThrow(json, schema)</code>	Parse or throw typed error
<code>parseEventBody(body, schema?)</code>	Parse API Gateway body
<code>parseJsonField(value, fallback?)</code>	Parse nested DB fields
<code>safeJsonStringify(value)</code>	Stringify with circular ref handling

## 62.7 Implementation Files

File	Purpose
<code>lambda/admin/code-quality.ts</code>	Admin API handler
<code>lambda/shared/utils/safe-json.ts</code>	Safe JSON utilities
<code>migrations/V2026_01_20_002__code_quality_metrics.sql</code>	Database schema
<code>apps/admin-dashboard/app/(dashboard)/code-quality/page.tsx</code>	Radiant dashboard
<code>apps/admin-dashboard/app/(dashboard)/thinktank/code-quality/page.tsx</code>	Think Tank dashboard

## Section 63: The Sovereign Mesh (PROMPT-36)

### 63.1 Overview

**”Every Node Thinks. Every Connection Learns. Every Workflow Assembles Itself.”**

The Sovereign Mesh introduces parametric AI assistance at every node level. Each Method, Agent, Service, and Connector can independently use AI to:

- **Disambiguate** unclear inputs
- **Infer** missing parameters
- **Recover** from errors intelligently
- **Validate** before execution
- **Explain** what was done

### 63.2 The Four Node Types

Type	Purpose	AI Capability
<b>Methods</b>	Deterministic reasoning, judging, transforming	Disambiguate, explain decisions
<b>Agents</b>	Goal-oriented work with OODA loops	Full reasoning, sub-task spawning
<b>Services</b>	Actions in external systems (3,000+ apps)	Infer params, recover errors, validate
<b>Libraries</b>	Local code execution	Generate code, explain transformations

### 63.3 Agent Registry

**Location:** `/sovereign-mesh/agents`

Built-in agents with OODA-loop execution:

Agent	Category	Description	HITL
Research Agent	research	Web research and information synthesis	No
Coding Agent	coding	Code writing, debugging, refactoring	No
Data Analysis Agent	data	Dataset analysis and visualization	No

Agent	Category	Description	HITL
Lead Generation Agent	outreach	Prospect research and list building	Yes
Editor Agent	creative	Content review and improvement	No
Automation Agent	operations	Multi-step workflow execution	No

### 63.4 AI Helper Service

Parametric configuration for each component:

```
interface AIHelperConfig {
  enabled: boolean;
  disambiguation: { enabled: boolean; model: string; confidenceThreshold: number };
  parameterInference: { enabled: boolean; model: string; examples: array };
  errorRecovery: { enabled: boolean; model: string; maxAttempts: number };
  validation: { enabled: boolean; model: string; checks: array };
  explanation: { enabled: boolean; model: string };
}
```

### 63.5 App Registry

3,000+ app integrations from Activepieces and n8n with AI enhancement layer.

**Daily Sync:** Definitions synced at 2 AM UTC **Health Checks:** Top 100 apps checked hourly

### 63.6 HITL Approval Queues

Human-in-the-loop approval for high-stakes decisions:

- **Agent Plan Approval:** Review agent execution plans
- **High Cost Approval:** Operations exceeding cost thresholds
- **Safety Review:** Operations flagged by safety evaluation

**SLA Monitoring:** Every minute with automatic escalation

### 63.7 Transparency Layer

Complete Cato decision visibility:

- Decision events with reasoning chains
- War Room deliberation capture (phase-by-phase)
- Pre-computed explanations (summary, standard, detailed, audit tiers)
- Decision feedback and learned patterns

### 63.8 Execution History & Replay

Time-travel debugging capabilities:

- Step-by-step state snapshots
- Replay sessions with modified inputs
- Execution diffs and divergence detection
- Bookmarks and annotations

### 63.9 API Endpoints

**Base Path:** /api/admin/sovereign-mesh

Endpoint	Method	Purpose
/dashboard	GET	Overview metrics
/agents	GET/POST	List/create agents
/agents/:id	GET/PUT/DELETE	Agent management
/executions	GET/POST	List/start executions
/executions/:id/cancel	POST	Cancel execution
/executions/:id/resume	POST	Resume paused execution
/apps	GET	List apps (3,000+)
/apps/:id	GET	App details
/apps/:id/ai-config	PUT	Update AI enhancements
/connections	GET	List OAuth connections
/decisions	GET	List Cato decisions
/decisions/:id/war-room	GET	War Room deliberations
/approvals	GET	List pending approvals
/approvals/:id/approve	POST	Approve request
/approvals/:id/reject	POST	Reject request
/ai-helper/config	GET/PUT	AI Helper configuration
/ai-helper/usage	GET	Usage statistics

### 63.10 Database Schema

Migrations (V2026\_01\_20\_003 through V2026\_01\_20\_010):

Table	Purpose
agents	Agent registry with AI helper config
agent_executions	OODA execution state
agent_iteration_logs	Detailed iteration tracking
apps	3,000+ app registry
app_connections	OAuth/API credentials
app_learned_inferences	AI learning loop
ai_helper_calls	Usage tracking
ai_helper_cache	Response caching
ai_helper_config	System/tenant configuration
workflow_blueprints	Pre-flight provisioning
capability_checks	Capability verification
cato_decision_events	Decision transparency
cato_war_room_deliberations	War Room capture
hitl_queue_configs	Approval queue setup
hitl_approval_requests	Approval requests
execution_snapshots	Time-travel debugging
replay_sessions	Replay/what-if analysis

### 63.11 Implementation Files

**Services:** | File | Purpose | |-----|-----| | lambda/shared/services/sovereign-mesh/ai-helper.service.ts | AI Helper Service | | lambda/shared/services/sovereign-mesh/agent-runtime.service.ts | Agent Runtime | | lambda/shared/services/sovereign-mesh/notification.service.ts | Email/Slack/Webhook notifications | | lambda/shared/services/sovereign-mesh/snapshot-capture.service.ts | Execution state snapshots | | lambda/shared/services/sovereign-mesh/index.ts | Service exports |

**Lambda Functions:** | File | Purpose | |-----|-----| | lambda/admin/sovereign-mesh.ts | Admin API handler | | lambda/scheduled/app-registry-sync.ts | Daily app sync (2 AM UTC) | |

lambda/scheduled/app-health-check.ts | Hourly health check for top 100 apps | | lambda/scheduled/hitl-sla-monitor.  
 | SLA monitoring with notifications | | lambda/workers/agent-execution-worker.ts | SQS-triggered  
 OODA processing | | lambda/workers/transparency-compiler.ts | Pre-compute decision explanations |

**CDK Infrastructure:** | File | Purpose | |-----|-----| | lib/stacks/sovereign-mesh-stack.ts | Com-  
 plete CDK stack with SQS queues, Lambdas, IAM |

**Dashboard Pages:** | File | Purpose | |-----|-----| | apps/admin-dashboard/app/(dashboard)/sovereign-mesh/page.tsx  
 | Main overview | | apps/admin-dashboard/app/(dashboard)/sovereign-mesh/agents/page.tsx | Agent  
 registry | | apps/admin-dashboard/app/(dashboard)/sovereign-mesh/apps/page.tsx | App browser | |  
 apps/admin-dashboard/app/(dashboard)/sovereign-mesh/transparency/page.tsx | Decision explorer  
 | | apps/admin-dashboard/app/(dashboard)/sovereign-mesh/ai-helper/page.tsx | AI Helper config |

**Database Migrations:** | File | Purpose | |-----|-----| | migrations/V2026\_01\_20\_003\_\_sovereign\_mesh\_agents.sql  
 | Agent schema | | migrations/V2026\_01\_20\_004\_\_sovereign\_mesh\_apps.sql | App schema | |  
 migrations/V2026\_01\_20\_005\_\_sovereign\_mesh\_ai\_helper.sql | AI Helper schema | | migrations/V2026\_01\_20\_006\_\_s  
 | Pre-flight provisioning | | migrations/V2026\_01\_20\_007\_\_sovereign\_mesh\_transparency.sql  
 | Transparency layer | | migrations/V2026\_01\_20\_008\_\_sovereign\_mesh\_hitl.sql | HITL ap-  
 proval queues | | migrations/V2026\_01\_20\_009\_\_sovereign\_mesh\_replay.sql | Execution replay |  
 | migrations/V2026\_01\_20\_010\_\_sovereign\_mesh\_seed.sql | Seed data

## Section 64: HITL Orchestration Enhancements

### 64.1 Overview

Advanced Human-in-the-Loop orchestration implementing industry best practices for intelligent question management.

**Philosophy:** "Ask only what matters. Batch for convenience. Never interrupt needlessly."

### 64.2 Core Features

Feature	Description
<b>SAGE-Agent Bayesian VOI</b>	Value-of-Information calculation to determine question necessity
<b>MCP Elicitation Schema</b>	Standardized question/response formats
<b>Question Batching</b>	Three-layer batching (time-window, correlation, semantic)
<b>Rate Limiting</b>	Global (50 RPM), per-user (10 RPM), per-workflow (5 RPM)
<b>Abstention Detection</b>	Output-based uncertainty detection for external models
<b>Question Deduplication</b>	TTL cache with fuzzy matching
<b>Escalation Chains</b>	Configurable multi-level escalation paths
<b>Two-Question Rule</b>	Max 2 clarifications per workflow

### 64.3 SAGE-Agent Bayesian VOI

The VOI service calculates whether asking a question provides enough expected value:

$VOI = Expected\_Information\_Gain - Ask\_Cost$

Decision =  $VOI > Threshold$  ? "ask" : "skip\_with\_default"

**Key Metrics:**

- Prior entropy calculation using Shannon entropy
- Expected posterior entropy estimation
- Ask cost based on urgency and workflow type
- Decision improvement impact weighting

## 64.4 Question Types (MCP Elicitation)

Type	Description
yes_no	Binary true/false question
single_choice	Select one from options
multiple_choice	Select multiple from options
free_text	Open-ended text response
numeric	Numeric value input
date	Date selection
confirmation	Explicit confirmation
structured	JSON schema-validated response

## 64.5 Abstention Detection Methods

For external models (no internal state access):

Method	Description
<b>Confidence Prompting</b>	Ask model to rate confidence 0-100
<b>Self-Consistency</b>	Sample N responses, measure agreement
<b>Semantic Entropy</b>	Cluster outputs, high entropy = uncertain
<b>Refusal Detection</b>	Detect hedging language patterns

**Future:** Linear probe abstention for self-hosted models via inference wrappers.

## 64.6 Rate Limiting Configuration

Scope	Requests/Min	Concurrent	Burst
Global	50	20	10
Per User	10	3	2
Per Workflow	5	2	1

## 64.7 Admin API Endpoints

Base: /api/admin/hitl-orchestration

Endpoint	Method	Description
/dashboard	GET	Complete dashboard data
/voi/statistics	GET	VOI decision statistics
/abstention/config	GET/PUT	Abstention settings
/abstention/statistics	GET	Abstention event stats
/batching/statistics	GET	Batch metrics
/rate-limits	GET	Rate limit configs
/rate-limits/:scope	PUT	Update rate limit
/escalation-chains	GET/POST	Manage escalation chains
/deduplication/statistics	GET	Cache statistics
/deduplication/invalidate	POST	Invalidate cache entries

64.8 Admin Dashboard

Location: /hitl-orchestration

Tabs:

- **Overview:** Key metrics, VOI breakdown, abstention reasons
- **Value of Information:** SAGE-Agent VOI statistics
- **Abstention Detection:** Detection methods, model statistics
- **Question Batching:** Batching strategies and metrics
- **Rate Limits:** Configuration table
- **Settings:** Threshold configuration

64.9 Database Tables

Table	Purpose
hitl_question_batches	Question batch records
hitl_rate_limits	Rate limit configuration
hitl_question_cache	Deduplication cache
hitl_voi_aspects	VOI aspect tracking
hitl_voi_decisions	VOI decision records
hitl_abstention_config	Abstention settings
hitl_abstention_events	Abstention event log
hitl_escalation_chains	Escalation chain configuration

64.10 Implementation Files

**Services:** | File | Purpose | |-----|-----| | lambda/shared/services/hitl-orchestration/mcp-elicitation.service.ts | Main orchestration | | lambda/shared/services/hitl-orchestration/voi.service.ts | Bayesian VOI | | lambda/shared/services/hitl-orchestration/abstention.service.ts | Uncertainty detection | | lambda/shared/services/hitl-orchestration/batching.service.ts | Question batching | | lambda/shared/services/hitl-orchestration/rate-limiting.service.ts | Rate limits | | lambda/shared/services/hitl-orchestration/deduplication.service.ts | Answer caching | | lambda/shared/services/hitl-orchestration/escalation.service.ts | Escalation chains |

**Lambda Functions:** | File | Purpose | |-----|-----| | lambda/admin/hitl-orchestration.ts | Admin API handler |

**Dashboard Pages:** | File | Purpose | |-----|-----| | apps/admin-dashboard/app/(dashboard)/hitl-orchestration/page | Radiant Admin | | apps/thinktank-admin/app/hitl-orchestration/page.tsx | Think Tank Admin |

**Database Migration:** | File | Purpose | |-----|-----| | migrations/V2026\_01\_20\_011\_\_hitl\_orchestration\_enhancement | Schema changes | | migrations/V2026\_01\_20\_012\_\_hitl\_semantic\_deduplication.sql | Semantic deduplication |

64.11 Semantic Deduplication (v5.34.0)

Enhanced question deduplication using pgvector embeddings for semantic similarity matching.

How It Works:

1. Questions are embedded using AI (1536-dimensional vectors)
2. Similar questions found via HNSW index cosine similarity search
3. Falls back to fuzzy matching if embeddings unavailable

Configuration:

Setting	Default	Description
<code>enableSemanticMatching</code>	<code>false</code>	Enable pgvector semantic search
<code>semanticSimilarityThreshold</code>	<code>0.85</code>	Minimum cosine similarity (0.0-1.0)
<code>maxSemanticCandidates</code>	<code>20</code>	Max candidates to check

#### API Endpoints:

Endpoint	Method	Description
<code>/semantic-deduplication/config</code>	GET	Get semantic config
<code>/semantic-deduplication/config</code>	PUT	Update semantic config
<code>/semantic-deduplication/stats</code>	GET	Match statistics (24h)
<code>/semantic-deduplication/backfill</code>	POST	Trigger embedding backfill

**Dashboard Tab:** "Deduplication" in HITL Orchestration admin

#### Match Statistics:

- Exact matches (hash-based)
- Fuzzy matches (Jaccard similarity)
- Semantic matches (pgvector cosine)
- Questions with embeddings count
- Average semantic similarity score

### 64.12 Scout HITL Integration (v5.34.0)

Bridges Cato's Scout persona (epistemic uncertainty mode) with HITL orchestration for intelligent clarification.

#### Flow:

1. Scout persona activates due to epistemic uncertainty
2. ScoutHITLIntegration generates prioritized clarification questions
3. Questions filtered through VOI scoring
4. High-VOI questions go to HITL, low-VOI get assumptions
5. Responses reduce uncertainty, allowing Scout to proceed

#### Domains:

- `medical` - HIPAA-sensitive, safety-critical
- `financial` - SOC2/PCI compliance
- `legal` - Regulatory compliance
- `bioinformatics` - Research accuracy
- `general` - Default domain

#### Aspect Impact Scores:

Aspect	Base Impact	Domain Boosts
<code>safety</code>	<code>0.95</code>	medical, bioinformatics
<code>compliance</code>	<code>0.90</code>	medical, financial, legal
<code>irreversible</code>	<code>0.85</code>	(all)
<code>cost</code>	<code>0.80</code>	financial
<code>accuracy</code>	<code>0.75</code>	medical, legal, bioinformatics
<code>timeline</code>	<code>0.60</code>	(none)

## API Endpoints (Cato Admin):

Base: /api/admin/cato

Endpoint	Method	Description
/scout-hitl/config	GET	Get Scout HITL config
/scout-hitl/config	PUT	Update config
/scout-hitl/sessions	GET	Recent clarification sessions
/scout-hitl/statistics	GET	Session statistics
/scout-hitl/domain-boosts	GET	Aspect domain boosts
/scout-hitl/domain-boosts	PUT	Update domain boosts

**Dashboard:** Cato → Scout HITL

### Configuration:

Setting	Default	Description
enabled	true	Enable Scout HITL integration
voiThreshold	0.3	Minimum VOI to ask question
maxQuestionsPerSession	3	Max clarifications before assuming
defaultDomain	general	Fallback domain

### Session Recommendations:

- **proceed** - Uncertainty resolved sufficiently
- **wait** - Still uncertain, user should wait
- **abort** - Critical uncertainty, cannot proceed safely

## 64.13 Flyte HITL Task Wrappers (v5.34.0)

Python task wrappers for easy HITL integration in Flyte workflows.

### Available Functions:

```
from radiant_flyte.utils import (
    ask_confirmation,
    ask_choice,
    ask_batch,
    ask_free_text,
)

# Confirmation (yes/no)
approved = await ask_confirmation(
    question="Deploy to production?",
    context={"environment": "prod"},
    timeout_seconds=300
)

# Single choice
selected = await ask_choice(
    question="Select deployment strategy",
    options=["rolling", "blue-green", "canary"],
    default="rolling"
)
```



```
# Batch questions
responses = await ask_batch([
    {"question": "Confirm rollback?", "type": "yes_no"},
    {"question": "Reason", "type": "free_text"},
])

# Free text
reason = await ask_free_text(
    question="Describe the issue",
    max_length=500
)
```

**Files:** | File | Purpose | |-----|-----| | packages/flyte/utils/hitl\_tasks.py | Task wrappers |

65. RAWS v1.1 - Model Selection System

RAWS (RADIANT AI Weighted Selection) provides intelligent real-time model selection using 8-dimension scoring across 13 weight profiles and 7 domains.

65.1 Overview

Component	Count	Description
Dimensions	8	Quality, Cost, Latency, Capability, Reliability, Compliance, Availability, Learning
Profiles	13	4 Optimization + 6 Domain + 3 SOFAI
Domains	7	Healthcare, Financial, Legal, Scientific, Creative, Engineering, General
Models	106+	50 external APIs + 56 self-hosted

65.2 Weight Profiles

Profile	Category	Primary Focus	Compliance
BALANCED	Optimization	Default, general purpose	-
QUALITY_FIRST	Optimization	Maximum accuracy	-
COST_OPTIMIZED	Optimization	Budget-conscious	-
LATENCY_CRITICAL	Optimization	Real-time applications	-
HEALTHCARE	Domain	Medical/clinical	HIPAA required
FINANCIAL	Domain	Finance/investment	SOC 2 required
LEGAL	Domain	Contracts/litigation	SOC 2 required
SCIENTIFIC	Domain	Research/academic	Optional
CREATIVE	Domain	Content/marketing	None
ENGINEERING	Domain	Code/software	Optional
SYSTEM_1	SOFAI	Fast, simple queries	-
SYSTEM_2	SOFAI	Complex reasoning	-
SYSTEM_2_5	SOFAI	Maximum reasoning	-

65.3 Domain Compliance Matrix

Domain	Required	Optional	Truth Engine	ECD Threshold
healthcare	HIPAA	FDA 21 CFR Part 11	Required	0.05

Domain	Required	Optional	Truth Engine	ECD Threshold
financial	SOC 2 Type II	PCI-DSS, GDPR, SOX	Required	0.05
legal	SOC 2 Type II	GDPR, State Bar	Required	0.05
scientific	None	FDA 21 CFR, GLP, IRB	Optional	0.08
creative	None	FTC Guidelines	Not Required	0.20
engineering	None	SOC 2, ISO 27001, NIST	Optional	0.10
general	None	None	Not Required	0.10

## 65.4 Admin API Endpoints

Base: /api/admin/raws

Endpoint	Method	Description
/select	POST	Select optimal model
/profiles	GET	List all 13 weight profiles
/profiles	POST	Create custom profile
/models	GET	List available models
/domains	GET	List 7 domain configurations
/detect-domain	POST	Test domain detection
/health	GET	Provider health status
/audit	GET	Selection audit log

## 65.5 CLI Commands

### # Profiles

```
radiant-cli raws profiles list --env production
radiant-cli raws profiles get HEALTHCARE --env production
```

### # Domains

```
radiant-cli raws domains list --env production
radiant-cli raws domains get healthcare --env production
```

### # Compliance

```
radiant-cli raws compliance summary --env production
radiant-cli raws models list --compliance HIPAA --env production
```

### # Audit

```
radiant-cli raws audit search --domain healthcare --last 24h --env production
```

## 65.6 Detailed Documentation

Document	Purpose
RAWS-ENGINEERING.md	Technical reference for engineers
RAWS-ADMIN-GUIDE.md	Operations and compliance guide
RAWS-USER-GUIDE.md	API guide for developers

## 65.7 Key Files

File	Purpose
migrations/V2026_01_21_004__raws_weighted_selection.sql	Database schema
lambda/shared/services/raws/types.ts	TypeScript types
lambda/shared/services/raws/domain-detector.service.ts	Domain detection
lambda/shared/services/raws/weight-profile.service.ts	Profile management
lambda/shared/services/raws/selection.service.ts	Main selection logic
lambda/admin/raws.ts	Admin API handler

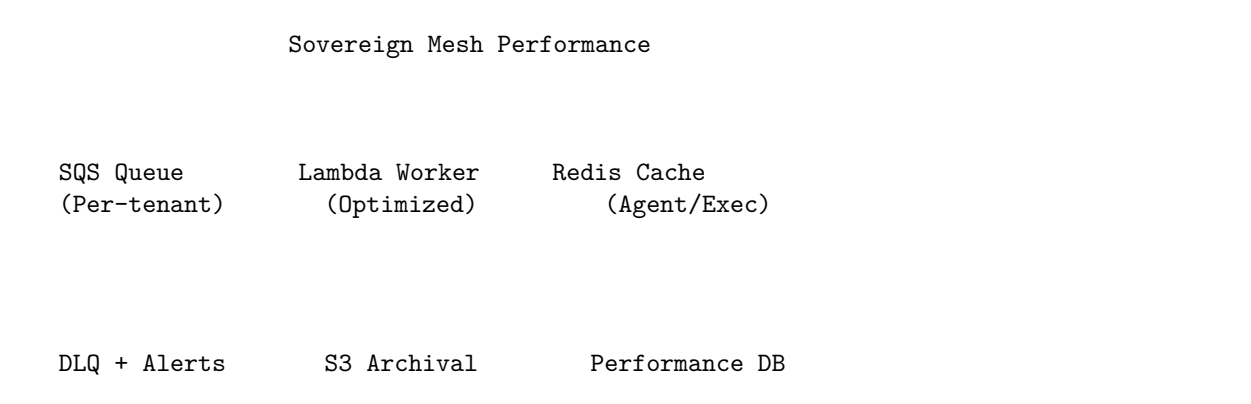
## 66. Sovereign Mesh Performance Optimization

The Sovereign Mesh Performance system provides comprehensive monitoring, configuration, and optimization for autonomous agent execution at scale. It enables administrators to tune Lambda concurrency, caching strategies, tenant isolation, and alert thresholds.

### 66.1 Overview

Component	Description
<b>SQS Dispatcher</b>	Actual message dispatch for OODA loop iterations
<b>Redis Cache</b>	Agent/execution state caching for reduced DB load
<b>Performance Config</b>	Per-tenant configuration with defaults
<b>Artifact Archival</b>	S3/hybrid storage for completed execution artifacts
<b>Rate Limiting</b>	Per-tenant and per-user concurrency limits

### 66.2 Architecture



### 66.3 Admin Dashboard

Navigate to **Sovereign Mesh** → **Performance** in the Admin Dashboard.

#### Overview Tab

- **Health Score:** 0-100 composite health metric
- **Active/Pending Executions:** Real-time counts
- **Queue Metrics:** Messages pending, in-flight, DLQ
- **Cache Hit Rate:** Redis/memory cache performance

- **OODA Phase Timing:** Per-phase latency breakdown
- **Cost Estimate:** Monthly Lambda/SQS projections

### Scaling Tab

Setting	Default	Range	Description
Max Concurrency	50	1-200	SQS event source concurrency
Provisioned Concurrency	5	0-50	Pre-warmed Lambda instances
Memory (MB)	2048	512-4096	Lambda memory allocation
Isolation Mode	shared	shared/dedicated/fifo	Tenant queue isolation
Max Per Tenant	50	1-100	Concurrent executions per tenant
Max Per User	10	1-25	Concurrent executions per user

### Caching Tab

- **Backend:** memory (Lambda-local) or redis (ElastiCache)
- **Hit Rate:** Target >80% for optimal performance
- **Cache Actions:** Clear tenant cache, view statistics

### Alerts Tab

Alert Type	Default Threshold	Description
DLQ Threshold	10 messages	Alert when DLQ exceeds threshold
Latency Threshold	30 seconds	Alert on slow executions
Budget Threshold	80%	Alert on budget exhaustion

**Recommendations Tab** AI-generated performance recommendations with one-click apply:

- Increase concurrency when utilization is high
- Improve cache hit rate suggestions
- Memory optimization recommendations

## 66.4 Configuration API

Base: `/api/admin/sovereign-mesh/performance`

Endpoint	Method	Description
<code>/dashboard</code>	GET	Get complete dashboard data
<code>/config</code>	GET	Get current configuration
<code>/config</code>	PUT/PATCH	Update configuration
<code>/recommendations</code>	GET	Get AI recommendations
<code>/recommendations/:id/apply</code>	POST	Apply a recommendation
<code>/alerts</code>	GET	List active alerts
<code>/alerts/:id/acknowledge</code>	POST	Acknowledge alert
<code>/alerts/:id/resolve</code>	POST	Resolve alert
<code>/cache/stats</code>	GET	Get cache statistics
<code>/cache</code>	DELETE	Clear tenant cache
<code>/queue/metrics</code>	GET	Get queue metrics
<code>/health</code>	GET	Health check

## 66.5 Database Tables

Table	Purpose
sovereign_mesh_performance_config	Per-tenant performance settings
sovereign_mesh_performance_alerts	Active and historical alerts
sovereign_mesh_performance_metrics	Time-series metrics (BRIN indexed)
sovereign_mesh_artifact_archives	Archived execution artifacts
sovereign_mesh_tenant_queues	Dedicated tenant queue mappings
sovereign_mesh_rate_limits	Rate limiting window counters
sovereign_mesh_config_history	Configuration change audit trail

## 66.6 Performance Indexes

The migration adds optimized indexes for common query patterns:

```
-- Fast tenant+status queries
CREATE INDEX idx_agent_executions_tenant_status ON agent_executions(tenant_id, status);

-- Fast agent+status queries
CREATE INDEX idx_agent_executions_agent_status ON agent_executions(agent_id, status);

-- Time-based queries
CREATE INDEX idx_agent_executions_created_at ON agent_executions(created_at DESC);

-- Partial index for running executions only
CREATE INDEX idx_agent_executions_running ON agent_executions(tenant_id, started_at)
WHERE status = 'running';

-- BRIN index for time-series metrics
CREATE INDEX idx_perf_metrics_tenant_time ON sovereign_mesh_performance_metrics
USING BRIN (tenant_id, metric_time);
```

## 66.7 Lambda Configuration

Production-optimized Lambda settings:

Setting	Value	Rationale
Memory	2048 MB	Complex OODA loop processing
Timeout	15 minutes	Long-running agent executions
Reserved Concurrency	100	Guaranteed capacity
Provisioned Concurrency	5	Eliminate cold starts
SQS Batch Size	1	OODA loop atomicity
SQS Max Concurrency	50	High throughput

## 66.8 Artifact Archival

Configure artifact archival for cost optimization:

Setting	Default	Description
Storage Backend	hybrid	database (small), s3 (large)
Archive After Days	7	Days until archival
Delete After Days	90	Days until deletion (0=never)

Setting	Default	Description
Max DB Bytes	65536	Threshold for S3 (64KB)
Compression	gzip	Compression algorithm

## 66.9 Rate Limiting

Built-in rate limiting protects against runaway executions:

```
-- Check if execution allowed
SELECT * FROM can_start_execution(:tenant_id, :user_id);

-- Returns: allowed, reason, current_count, max_allowed
```

## 66.10 Key Files

File	Purpose
packages/shared/src/types/sovereign-mesh-performance.types.ts	TypeScript types
packages/infrastructure/migrations/V2026_01_21_001__sovereign_mesh_performance.sql	Database schema
lambda/shared/services/sovereign-mesh/sqs-dispatcher.service.ts	SQS message dispatch
lambda/shared/services/sovereign-mesh/redis-cache.service.ts	Redis/memory caching
lambda/shared/services/sovereign-mesh/performance-config.service.ts	Configuration manager
lambda/shared/services/sovereign-mesh/artifact-archival.service.ts	S3 archival
lambda/admin/sovereign-mesh-performance.ts	Admin API handler
apps/admin-dashboard/app/(dashboard)/sovereign-mesh/performance/page.tsx	Admin UI

## 66.11 Troubleshooting

Issue	Cause	Solution
High DLQ count	Processing failures	Check CloudWatch logs, increase timeout
Low cache hit rate	Cold starts, small TTL	Increase TTL, enable cache warming
Slow executions	Memory constraints	Increase Lambda memory
Rate limit errors	Too many concurrent	Increase per-tenant limits
Queue backlog	Insufficient concurrency	Increase maxConcurrency

## 67. Infrastructure Scaling (100 to 500K Sessions)

Comprehensive infrastructure scaling system enabling seamless growth from development (100 sessions) to enterprise scale (500,000+ concurrent sessions) with real-time cost visibility.

### 67.1 Scaling Tiers

Tier	Target Sessions	Monthly Cost	Use Case
<b>Development</b>	100	~\$70	Testing, development, POC
<b>Staging</b>	1,000	~\$500	Pre-production testing
<b>Production</b>	10,000	~\$5,000	Standard production workloads
<b>Enterprise</b>	500,000	~\$68,500	Global scale, multi-region

## 67.2 Admin Dashboard

Navigate to **Sovereign Mesh** → **Scaling** in the Admin Dashboard.

### Overview Tab

- **Active Sessions:** Real-time count with utilization gauge
- **Peak Sessions:** Today, week, and month peaks
- **Bottleneck Indicator:** Current limiting component
- **Cost per Session:** Real-time cost efficiency metric
- **Component Health:** Status of Lambda, Aurora, Redis, API Gateway, SQS

### Sessions Tab

- **Session Capacity:** Current vs maximum with headroom
- **Session Statistics:** Historical counts and averages
- **Capacity by Component:** Per-component session limits

### Infrastructure Tab

- **Lambda Configuration:** Reserved/provisioned concurrency, memory, timeout
- **Aurora Configuration:** ACU range, read replicas, global database
- **Redis Configuration:** Node type, shards, cluster mode
- **API Gateway Configuration:** Rate limits, CloudFront

### Cost Tab

- **Cost Breakdown:** Per-component monthly costs with progress bars
- **Cost Metrics:** Cost per session, per 1000 sessions, annual estimate
- **Component Costs:** Lambda, Aurora, Redis, API Gateway, SQS, CloudFront, Data Transfer

### Scale Tab

- **Quick Scale:** One-click tier selection with instant apply
- **Scaling Comparison:** Cost delta between tiers
- **What Changes:** Detailed list of infrastructure modifications

## 67.3 Scaling Profiles

Each scaling tier configures multiple components:

### Development (Scale-to-Zero)

Lambda:	0 provisioned, 10 max concurrent, 1024 MB
Aurora:	0.5-2 ACU, 0 replicas, no global
Redis:	cache.t4g.micro, 1 shard, no cluster
API:	100 RPS, no CloudFront
SQS:	2 standard queues

### Staging

Lambda:	0 provisioned, 50 max concurrent, 2048 MB
Aurora:	1-8 ACU, 1 replica, no global
Redis:	cache.t4g.small, 1 shard, 1 replica
API:	1,000 RPS, no CloudFront
SQS:	5 standard, 2 FIFO queues

## Production

Lambda: 5 provisioned, 200 max concurrent, 2048 MB  
Aurora: 4-64 ACU, 2 replicas, PgBouncer  
Redis: cache.r6g.large, 1 shard, 2 replicas  
API: 10,000 RPS, CloudFront enabled  
SQS: 10 standard, 5 FIFO queues

## Enterprise (500K)

Lambda: 100 provisioned, 1000 max concurrent, 3072 MB  
Aurora: 16-256 ACU, 3 replicas, Global Database (5 regions)  
Redis: cache.r6g.xlarge, 10 shards, 2 replicas, cluster mode  
API: 100,000 RPS, CloudFront, 5 regional endpoints  
SQS: 50 standard, 50 FIFO queues

### 67.4 Cost Calculation

Real-time cost estimation based on AWS pricing:

Component	Pricing Model	Example (Production)
<b>Lambda Provisioned</b>	\$0.000004167/GB-second	$5 \times 2\text{GB} \times 24\text{h} \times 30\text{d} = \sim\$360/\text{mo}$
<b>Aurora ACU</b>	\$0.12/ACU-hour	$34 \text{ avg ACU} \times 24\text{h} \times 30\text{d} = \sim\$2,900/\text{mo}$
<b>Redis</b>	\$0.182/hour (r6g.large)	$3 \text{ nodes} \times 24\text{h} \times 30\text{d} = \sim\$390/\text{mo}$
<b>API Gateway</b>	\$1.00/million requests	$10\% \text{ of } 10\text{K RPS} = \sim\$260/\text{mo}$
<b>SQS</b>	\$0.40/million standard	$15 \text{ queues} \times 1\text{M} = \sim\$6/\text{mo}$
<b>CloudFront</b>	\$0.085/GB + \$0.01/10K req	$\sim\$500/\text{mo}$
<b>Data Transfer</b>	\$0.02/GB inter-region	$\sim\$200/\text{mo (global)}$

### 67.5 Session Capacity Calculation

Maximum concurrent sessions is limited by the bottleneck component:

Lambda Max = `maxConcurrency` × 10 sessions/concurrent  
Aurora Max = `connectionPoolSize` × 50 sessions/connection  
Redis Max = `maxConnections`  
API Gateway Max = `throttlingRateLimit`

Effective Max = MIN(Lambda, Aurora, Redis, API Gateway)

### 67.6 Auto-Scaling Rules

Configure automatic scaling based on metrics:

Metric	Condition	Action
<code>session_count</code> > 80% capacity	Duration: 5 min	Scale up
<code>session_count</code> < 20% capacity	Duration: 30 min	Scale down
<code>cpu_utilization</code> > 70%	Duration: 5 min	Increase concurrency
<code>latency_p99</code> > 5000ms	Duration: 2 min	Increase memory
<code>queue_depth</code> > 1000	Duration: 1 min	Scale out workers

### 67.7 Scheduled Scaling

Pre-configure scaling for predictable patterns:



```
# Scale up for business hours (PST)
0 6 * * MON-FRI → Production tier
# Scale down after hours
0 18 * * MON-FRI → Staging tier
# Minimal on weekends
0 0 * * SAT → Development tier
```

## 67.8 API Endpoints

Base: /api/admin/sovereign-mesh/scaling

Endpoint	Method	Description
/dashboard	GET	Complete scaling dashboard
/profiles	GET	List all profiles
/profiles/active	GET	Get active profile
/profiles	POST	Create new profile
/profiles/:id	PUT	Update profile
/profiles/:id/apply	POST	Apply profile
/sessions	GET	Session metrics
/sessions/capacity	GET	Capacity info
/sessions/trends	GET	Historical trends
/cost	GET	Current cost estimate
/cost/estimate	POST	Estimate custom config
/operations	GET	Recent operations
/alerts	GET	Active alerts
/health	GET	Component health
/presets	GET	Available presets
/presets/:tier/apply	POST	Apply preset tier

### 67.8.1 Cost Analytics Integration

Infrastructure scaling costs are integrated into the AWS Cost Monitoring service and appear as a **line item group** in the Cost Analytics page.

**Endpoint:** GET /api/admin/costs/infrastructure-scaling

**Response:**

```
{
  "scalingCosts": {
    "groupName": "Infrastructure Scaling",
    "totalCost": 5234.56,
    "estimatedMonthlyCost": 5000.00,
    "tier": "production",
    "targetSessions": 10000,
    "costPerSession": 0.5234,
    "lineItems": [
      {
        "component": "Aurora",
        "description": "Serverless v2 (4-64 ACU, 2 replicas)",
        "baseCost": 2916.00,
        "usageCost": 0,
        "totalCost": 2916.00,
        "unit": "ACU-hour",
        "quantity": 24336,

```

```

        "pricePerUnit": 0.12,
        "percentage": 55.7
    },
    // ... other components
],
"lastUpdated": "2026-01-21T19:30:00Z"
}
}

```

**Line Item Components:** | Component | Description | Pricing Model | |-----|-----|-----|  
| **Lambda** | Provisioned concurrency | \$0.000004167/GB-second | | **Aurora** | Serverless v2 ACU-hours |  
\$0.12/ACU-hour | | **Redis** | ElastiCache node-hours | Varies by node type | | **API Gateway** | HTTP  
API requests | \$1.00/million requests | | **SQS** | Queue requests | \$0.40-0.50/million | | **CloudFront** | Edge  
distribution | \$0.085/GB + requests |

**Cost Analytics UI:** Navigate to **Costs** in the Admin Dashboard to see:

- Infrastructure Scaling section with tier badge
- Summary cards: Total Monthly, Estimated, Cost/Session, Components
- Detailed line items table with percentage bars
- Per-component breakdown with icons

## 67.9 Database Tables

Table	Purpose
sovereign_mesh_scaling_profiles	Scaling profile configurations
sovereign_mesh_session_metrics	Real-time session metrics (1-min)
sovereign_mesh_session_metrics_hourly	Aggregated hourly metrics
sovereign_mesh_scaling_operations	Operation history
sovereign_mesh_autoscaling_rules	Auto-scaling rules
sovereign_mesh_scheduled_scaling	Scheduled scaling events
sovereign_mesh_component_health	Component health snapshots
sovereign_mesh_scaling_alerts	Scaling alerts
sovereign_mesh_cost_records	Daily cost records

## 67.10 Key Files

File	Purpose
packages/shared/src/types/sovereign-mesh-scaling.types.ts	TypeScript types
packages/infrastructure/migrations/V2026_01_21_002__sovereign_mesh_scaling.sql	Database schema
lambda/shared/services/sovereign-mesh/scaling.service.ts	Scaling service
lambda/admin/sovereign-mesh-scaling.ts	Admin API handler
apps/admin-dashboard/app/(dashboard)/sovereign-mesh/scaling/page.tsx	Admin UI

## 67.11 Scaling Operations Workflow

1. Select Target Tier/Profile
2. Calculate Changes
  - Component diffs
  - Cost impact
  - Downtime assessment

3. Create Operation Record
  - Status: pending/in\_progress
  - Estimated duration
  - Rollback plan
4. Apply Changes (if approved)
  - Update Lambda concurrency
  - Modify Aurora ACUs
  - Resize Redis cluster
  - Update API Gateway limits
5. Verify & Complete
  - Health checks pass
  - Session capacity confirmed
  - Operation: completed

## 67.12 Troubleshooting

Issue	Cause	Solution
Scale operation stuck	CDK deployment in progress	Wait or rollback
Sessions exceeding capacity	Traffic spike	Increase tier or concurrency
High cost per session	Over-provisioned	Scale down during low traffic
Cold starts in production	No provisioned concurrency	Enable provisioned concurrency
Cross-region latency	Single region	Enable global database + Redis

## Section 68: Schema-Adaptive Reports (v5.39.0)

**Location:** Admin Dashboard → Reports → Schema Builder

Dynamic report builder that automatically discovers and adapts to database schema changes.

### 68.1 Overview

The Schema-Adaptive Report Writer provides:

- **Automatic Schema Discovery** - Queries `information_schema` to discover tables and columns
- **Intelligent Categorization** - Groups tables by domain (Core, AI, Billing, Analytics, System)
- **Dynamic Query Building** - Constructs SQL queries based on user selections
- **Visual Filter Builder** - 11 operators (=, >, <, LIKE, IN, BETWEEN, IS NULL, IS NOT NULL)
- **Date Presets** - Quick filters (Today, Yesterday, Last 7/30 Days, This/Last Month)
- **Per-Field Aggregation** - COUNT, SUM, AVG, MIN, MAX, COUNT DISTINCT per column
- **Sort & Group Builders** - Visual ORDER BY and GROUP BY configuration
- **SQL Preview** - Live-generated SQL query with dark-themed display
- **AI Suggestions** - Recommends useful report templates based on schema analysis
- **Visualization Toggles** - Table, Bar, Line, Pie chart view switches
- **Multi-tenant Security** - All queries respect RLS policies

## 68.2 AI Report Writer (v5.42.0)

Enterprise-grade AI-powered report generation with text and voice input, interactive charts, smart insights, and brand customization.

### Core Features:

- **Natural Language Generation** - Describe reports in plain English
- **Voice Input** - Web Speech API for hands-free report creation
- **AI Modification** - Refine reports with follow-up prompts
- **Report Styles** - Executive Summary, Detailed Analysis, Dashboard View, Narrative
- **Rich Formatting** - Headings, metrics cards, charts, tables, lists, quotes
- **Edit Mode** - Click sections to select, use format panel for styling
- **Undo/Redo** - Full history navigation
- **Export** - PDF, Excel, HTML, Print

### Interactive Charts (v5.42.0):

- Real Recharts visualizations (not placeholders)
- Bar, Line, Pie, Area chart types
- Auto-formatted tooltips (K/M suffixes)
- 8-color palette for data series
- Responsive container adapts to panel width

### Smart Insights (v5.42.0):

- AI-powered anomaly detection
- Trend analysis with predictions
- Achievement highlighting
- Actionable recommendations
- Warning alerts for concerning metrics
- Severity levels (low/medium/high)
- Confidence scores per insight

### Brand Kit (v5.42.0):

- Logo upload (PNG, JPG)
- Company name and tagline
- Primary/Secondary/Accent color pickers
- Font selection (Inter, Georgia, Roboto, etc.)
- Quick color presets (blue/green/purple/amber/slate)
- Live preview card
- Reset to defaults

### Usage:

1. Navigate to Reports → AI Writer tab
2. Select report style (Executive, Detailed, Dashboard, Narrative)
3. Type or speak your report request
4. Review generated report in preview panel
5. Use modification prompt to refine ("Add security metrics section")
6. Toggle Edit Mode to click and modify sections
7. Use Format panel for styling (bold, italic, alignment)
8. Export to PDF/Excel/HTML or Print

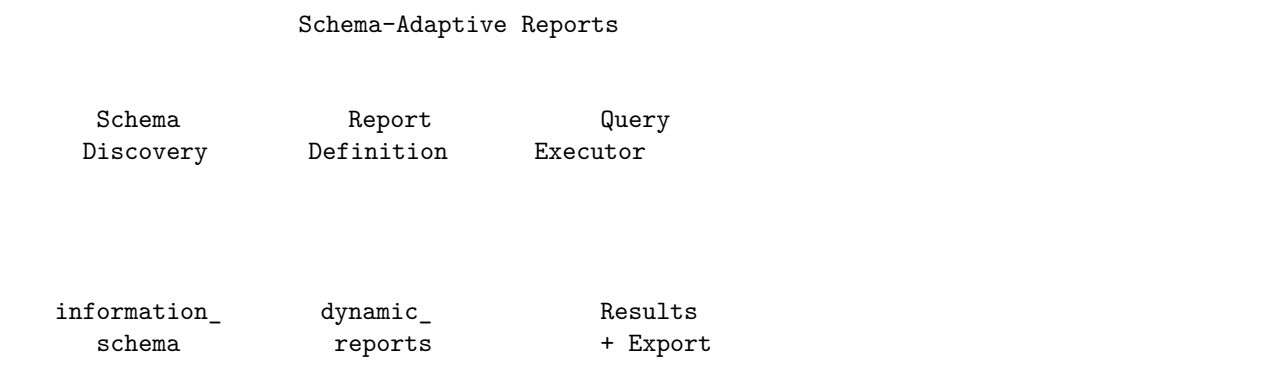
### Voice Commands:

- Click microphone icon to start voice input
- Speak naturally: "Create a monthly usage report with cost breakdown"
- Voice automatically transcribes to text input

- Press Enter or click Generate to process

**Report Sections:** | Type | Description | |-----|-----| | **heading** | H1-H3 headings with hierarchical styling | | **paragraph** | Body text with muted foreground | | **metrics** | 4-column KPI cards with trends | | **chart** | Placeholder for bar/line/pie/area charts | | **table** | Data tables with headers | | **list** | Bullet point lists | | **quote** | Blockquote with left border |

## 68.3 Architecture



## 68.3 Schema Discovery

The service discovers:

Element	Description
<b>Tables</b>	All user-accessible tables
<b>Columns</b>	Column names, types, nullability
<b>Primary Keys</b>	Identified for each table
<b>Foreign Keys</b>	Relationships between tables
<b>Indexes</b>	Available indexes for optimization

## 68.4 Table Categories

Category	Tables Included
<b>Core</b>	tenants, users, sessions, api_keys
<b>AI</b>	models, prompts, responses, brain_plans
<b>Billing</b>	subscriptions, credits, invoices, payments
<b>Analytics</b>	events, metrics, usage_logs
<b>System</b>	configurations, audit_logs, health_checks

## 68.5 Report Definition

```

interface DynamicReportDefinition {
    id?: string;
    name: string;
    description?: string;
    baseTable: string;
    fields: ReportField[];
    filters?: ReportFilter[];
}
  
```

```

joins?: ReportJoin[];
groupBy?: string[];
orderBy?: { column: string; direction: 'asc' | 'desc' }[];
limit?: number;
}

interface ReportField {
  column: string;
  alias?: string;
  aggregation?: 'count' | 'sum' | 'avg' | 'min' | 'max' | 'distinct';
  format?: 'text' | 'number' | 'currency' | 'percentage' | 'date' | 'datetime';
}

```

## 68.6 API Endpoints

Base: /api/admin/dynamic-reports

Method	Endpoint	Description
GET	/schema	Discover database schema with categorization
GET	/suggestions	AI-generated report suggestions
GET	/	List saved report definitions
POST	/	Save a new report definition
POST	/execute	Execute a report and return results
POST	/export	Export report results as CSV
DELETE	/:id	Delete a saved report

## 68.7 Request/Response Examples

### Schema Discovery:

GET /api/admin/dynamic-reports/schema

#### Response:

```

{
  "schema": [
    {
      "category": "Core",
      "tables": [
        {
          "name": "users",
          "columns": [
            { "name": "id", "type": "uuid", "nullable": false, "isPrimaryKey": true },
            { "name": "email", "type": "text", "nullable": false },
            { "name": "created_at", "type": "timestampz", "nullable": false }
          ]
        }
      ]
    }
  ]
}

```

### Execute Report:

POST /api/admin/dynamic-reports/execute

```

{

```

```

"baseTable": "users",
"fields": [
  { "column": "created_at", "format": "date" },
  { "column": "id", "aggregation": "count", "alias": "user_count" }
],
"groupBy": ["DATE(created_at)"],
"orderBy": [{ "column": "created_at", "direction": "desc" }],
"limit": 30
}

```

#### Response:

```

{
  "results": [
    { "created_at": "2026-01-21", "user_count": 145 },
    { "created_at": "2026-01-20", "user_count": 132 }
  ],
  "rowCount": 30,
  "executionTime": 45
}

```

## 68.8 Database Tables

```

-- Report definitions
CREATE TABLE dynamic_reports (
  id UUID PRIMARY KEY DEFAULT gen_random_uuid(),
  tenant_id UUID NOT NULL REFERENCES tenants(id),
  name TEXT NOT NULL,
  description TEXT,
  definition JSONB NOT NULL,
  created_by UUID REFERENCES users(id),
  created_at TIMESTAMPTZ DEFAULT NOW(),
  updated_at TIMESTAMPTZ DEFAULT NOW(),
  last_run_at TIMESTAMPTZ,
  run_count INTEGER DEFAULT 0
);

```

```

-- Execution history
CREATE TABLE dynamic_report_executions (
  id UUID PRIMARY KEY DEFAULT gen_random_uuid(),
  report_id UUID REFERENCES dynamic_reports(id),
  tenant_id UUID NOT NULL,
  executed_by UUID REFERENCES users(id),
  executed_at TIMESTAMPTZ DEFAULT NOW(),
  execution_time_ms INTEGER,
  row_count INTEGER,
  status TEXT DEFAULT 'completed',
  error_message TEXT
);

```

```

-- Scheduled reports
CREATE TABLE dynamic_report_schedules (
  id UUID PRIMARY KEY DEFAULT gen_random_uuid(),
  report_id UUID REFERENCES dynamic_reports(id),
  tenant_id UUID NOT NULL,

```

```

    schedule_cron TEXT NOT NULL,
    recipients JSONB,
    format TEXT DEFAULT 'csv',
    enabled BOOLEAN DEFAULT true,
    last_sent_at TIMESTAMPTZ,
    next_run_at TIMESTAMPTZ
);

```

## 68.9 Implementation Files

File	Purpose
packages/infrastructure/lambda/shared/services/schema-adaptive-reports.service.ts	Core service
packages/infrastructure/lambda/admin/dynamic-reports.ts	API handler
packages/infrastructure/migrations/V2026_01_21_003__dynamic_reports.sql	Database schema
apps/admin-dashboard/app/(dashboard)/reports/page.tsx	Admin UI with Schema

## 68.10 Security Considerations

- All queries run with tenant RLS context
- Schema discovery excludes system tables
- Parameterized queries prevent SQL injection
- Query timeout limits prevent resource exhaustion
- Audit logging for all report executions

## Section 69: Cato Genesis (v5.39.0)

**Location:** Admin Dashboard → Cato → Genesis

Autonomous AI genesis and self-improvement configuration for the Cato consciousness system.

### 69.1 Overview

Cato Genesis enables:

- **Autonomous Decision-Making** - AI makes decisions within configured boundaries
- **Self-Improvement** - System optimizes its own behavior over time
- **Safety Guardrails** - Configurable thresholds and restrictions
- **Human Oversight** - Approval requirements for critical actions

### 69.2 Configuration

Setting	Default	Description
enabled	true	Enable Genesis engine
autonomyLevel	65%	Decision-making independence (0-100%)
safetyThreshold	85%	Confidence required for autonomous action
learningRate	0.7	How aggressively to learn from outcomes
explorationEnabled	true	Allow exploration of new strategies
selfImprovementEnabled	false	Allow self-optimization (requires approval)
humanOversightRequired	true	Require human approval for critical actions
maxActionsPerMinute	100	Rate limiting for autonomous actions



### 69.3 Autonomy Levels

Level	Range	Description
<b>Restricted</b>	0-30%	Most actions require approval
<b>Monitored</b>	31-70%	Routine actions autonomous, critical reviewed
<b>Autonomous</b>	71-100%	Full autonomy within safety bounds

### 69.4 Safety Guardrails

#### Allowed Actions:

- Read operations
- Analysis tasks
- Recommendations
- Non-destructive updates

#### Restricted Actions (require approval):

- Delete operations
- External API calls
- User data modifications
- System configuration changes

### 69.5 Metrics

Metric	Description
<b>Total Decisions</b>	Cumulative decisions made
<b>Autonomous Actions</b>	Actions taken without human input
<b>Human Interventions</b>	Times humans overrode decisions
<b>Safety Violations</b>	Actions blocked by guardrails
<b>Learning Cycles</b>	Self-improvement iterations
<b>Avg Confidence</b>	Mean decision confidence score

### 69.6 Dashboard Tabs

Tab	Purpose
<b>Configuration</b>	Enable/disable, basic settings
<b>Autonomy Controls</b>	Autonomy level, rate limits
<b>Safety Guardrails</b>	Threshold configuration, action restrictions
<b>Learning &amp; Evolution</b>	Self-improvement settings, learning statistics

### 69.7 Implementation Files

File	Purpose
apps/admin-dashboard/app/(dashboard)/cato/genesis/page.tsx	Admin UI
packages/infrastructure/lambda/shared/services/cato/genesis.service.ts	Genesis service
packages/infrastructure/lambda/admin/cato-genesis.ts	API handler

## Section 70: Cortex Memory System (v4.20.0)

**Location:** Admin Dashboard → Memory → Cortex

Enterprise-scale tiered memory architecture for AI context and knowledge management.

### 70.1 Overview

The Cortex Memory System provides three-tier storage replacing direct database access:

Tier	Technology	Latency	Retention	Use Case
<b>Hot</b>	Redis + DynamoDB	< 10ms	4 hours	Session context, Ghost Vectors
<b>Warm</b>	Neptune + pgvector	< 100ms	90 days	Knowledge graph, entity relationships
<b>Cold</b>	S3 + Iceberg	< 2s	7+ years	Historical archives, compliance records

### 70.2 Dashboard Features

**Overview Page** (/cortex):

- Tier health cards with status indicators
- Data flow metrics (promotions, archivals, retrievals)
- Active alerts with acknowledgment
- Zero-Copy mount manager
- Housekeeping task status

**Quick Links:**

- Graph Explorer (/cortex/graph)
- Conflicts (/cortex/conflicts)
- GDPR Erasure (/cortex/gdpr)
- Settings (/cortex/settings)

### 70.3 Hot Tier Configuration

Setting	Default	Description
hot_redis_cluster_mode	true	Enable Redis sharding
hot_shard_count	3	Number of shards
hot_replicas_per_shard	2	HA replicas
hot_default_ttl_seconds	14400	4 hours default TTL
hot_overflow_to_dynamodb	true	Large value overflow

### 70.4 Warm Tier Configuration

Setting	Default	Description
warm_neptune_mode	serverless	Serverless or provisioned
warm_retention_days	90	Days before archival
warm_graph_weight_percent	60	Graph weight in hybrid search
warm_vector_weight_percent	40	Vector weight in hybrid search

### 70.5 Graph-RAG Knowledge

The Warm tier uses hybrid Graph-RAG search:

Hybrid Score = (Vector Similarity × 40%) + (Graph Traversal × 60%)

**Node Types:** document, entity, concept, procedure (evergreen), fact (evergreen), golden\_qa (verified answers)

**Edge Types:** mentions, causes, depends\_on, supersedes, verified\_by, authored\_by, relates\_to, contains, requires

70.5.1 Golden Rules (Override System)

Administrators can create high-priority rules that supersede all other data:

Rule Type	Description	Use Case
force_override	Always use this answer	"Max pressure is 100 PSI"
ignore_source	Never use this source	"Ignore Manual v1"
prefer_source	Prioritize this source	"Prefer 2026 specs"
deprecate	Mark as obsolete	"Manual v2024 superseded"

**Chain of Custody:** Every fact includes who verified it, when, and a digital signature for audit trail.

70.6 Zero-Copy Mounts & Stub Nodes

**The Innovation:** Connect to tenant data lakes without moving data. Creates **Stub Nodes** in the graph.

Source	Description	Connection
Snowflake	Data Share connection	OAuth + Data Share
Databricks	Delta Lake / Unity Catalog	Service Principal
S3	Customer S3 bucket	Cross-account IAM
Azure	Data Lake Gen2	Managed Identity

Stub Node Mechanism:

- Scans external storage metadata
- Creates lightweight graph nodes pointing to external content
- Content fetched **only** when Graph Traversal determines it's needed

**Actions:** Add Mount, Rescan, Delete

70.7 Curator Entrance Exams

Verify AI-extracted knowledge through SME validation:

Workflow:

1. AI ingests documents and extracts facts
2. System generates quiz questions from extracted facts
3. SME takes exam: Verify or Correct
4. Corrections automatically create Golden Rules
5. Verified facts get Chain of Custody signature

API:

```
# Generate exam for a domain
POST /api/admin/cortex/v2/exams
{
  "domainId": "hydraulics",
  "domainPath": "Engineering > Hydraulics",
```

```

    "questionCount": 10,
    "passingScore": 80
}

# Complete exam and process corrections
POST /api/admin/cortex/v2/exams/{examId}/complete

```

## 70.8 Live Telemetry (MQTT/OPC UA)

Inject real-time sensor data into AI context:

Protocol	Use Case	Example
MQTT	IoT sensors	mqtt://broker:1883
OPC UA	Industrial PLCs	opc.tcp://plc:4840
Kafka	Event streams	kafka://cluster:9092
WebSocket	Real-time dashboards	wss://server/feed

### Configuration:

```

POST /api/admin/cortex/v2/telemetry/feeds
{
  "name": "pump_302_sensors",
  "protocol": "opc_ua",
  "endpoint": "opc.tcp://plc.factory.local:4840",
  "nodeIds": ["ns=2;s=Pump302.Pressure"],
  "contextInjection": true
}

```

## 70.9 Model Migration

One-click swap between AI models without losing Cortex knowledge.

See Model Migration section above for supported models and workflow.

## 70.10 Housekeeping (Twilight Dreaming)

Task	Frequency	Purpose
TTL Enforcement	Hourly	Expire Hot tier keys
Archive Promotion	Nightly	Move Warm → Cold
Deduplication	Nightly	Merge duplicate nodes
Graph Expansion	Weekly	Infer missing links
Conflict Resolution	Nightly	Flag contradictions
Iceberg Compaction	Nightly	Optimize Cold storage
Index Optimization	Weekly	Reindex vectors

**Manual Trigger:** Click play button next to any task

## 70.8 GDPR Erasure

GDPR Article 17 "Right to be Forgotten" cascade deletion:

Tier	SLA	Method
Hot	Immediate	Redis key deletion
Warm	24 hours	Node status → deleted
Cold	72 hours	Tombstone records

Create Request: POST /api/admin/cortex/gdpr/erasure

## 70.9 Monitoring Thresholds

**Hot Tier:** | Metric | Warning | Critical | |-----|-----|-----| | Memory Usage | > 70% | > 85% | | Cache Hit Rate | < 90% | < 80% | | p99 Latency | > 5ms | > 10ms |

**Warm Tier:** | Metric | Warning | Critical | |-----|-----|-----| | Neptune CPU | > 70% | > 90% | | Graph Nodes | > 50M | > 100M |

## 70.10 API Endpoints

Base: /api/admin/cortex

GET	/overview	Full dashboard data
GET	/config	Tier configuration
PUT	/config	Update configuration
GET	/health	Tier health status
POST	/health/check	Trigger health check
GET	/alerts	Active alerts
POST	/alerts/:id/acknowledge	Acknowledge alert
GET	/metrics	Data flow metrics
GET	/graph/stats	Node/edge counts
GET	/graph/explore	Search graph nodes
GET	/graph/conflicts	Unresolved conflicts
GET	/housekeeping/status	Task statuses
POST	/housekeeping/trigger	Run task manually
GET	/mounts	Zero-Copy mounts
POST	/mounts	Create mount
POST	/mounts/:id/rescan	Rescan mount
DELETE	/mounts/:id	Delete mount
GET	/gdpr/erasure	Erasure requests
POST	/gdpr/erasure	Create erasure request

## 70.11 Implementation Files

File	Purpose
packages/shared/src/types/cortex-memory.types.ts	Type definitions
packages/infrastructure/migrations/V2026_01_23_002__cortex_memory_system.sql	Database schema
packages/infrastructure/lambda/shared/services/cortex/tier-coordinator.service.ts	Orchestration
packages/infrastructure/lambda/admin/cortex.ts	Admin API
apps/admin-dashboard/app/(dashboard)/cortex/page.tsx	Dashboard UI

## 70.12 The Sovereign Cortex Moats

The Cortex Memory System creates six interlocking competitive moats that make customer departure operationally prohibitive:

Moat	Mechanism	Why Customers S
<b>Semantic Structure</b>	Knowledge Graph with explicit relationships vs competitor "buckets of text"	Moving millions o
<b>Chain of Custody</b>	Every critical fact digitally signed by SME via Entrance Exam	Losing audit trail
<b>Tribal Delta</b>	Golden Rules capture real-world exceptions ("15 days in Mexico City, not 30")	Encoded intuition
<b>Sovereignty</b>	Model-agnostic: Cortex is the asset, models are disposable CPUs	"Better model? F
<b>Entropy Reversal</b>	Twilight Dreaming deduplicates, resolves conflicts, compresses nightly	Competitors get s
<b>Mentorship Equity</b>	Gamified Curator Quiz creates psychological ownership	After 50 hours of

### Administrative Implications:

- **Retention:** Monitor time invested per domain via Entrance Exam analytics
- **Value Demonstration:** Export Golden Rules count and relationship density for QBRs
- **Migration Prevention:** Highlight Chain of Custody signatures that would be lost

### 70.13 Related Documentation

- CORTEX-MEMORY-ADMIN-GUIDE.md - Full operations guide
- CORTEX-ENGINEERING-GUIDE.md - Technical reference

## 71. Semantic Blackboard & Multi-Agent Orchestration

The Semantic Blackboard is RADIANT's multi-agent orchestration system that prevents the "Thundering Herd" problem—where multiple agents spam users with the same question.

### 71.1 Overview

Feature	Description
<b>Semantic Question Matching</b>	Vector similarity search using OpenAI ada-002 embeddings
<b>Answer Reuse</b>	Auto-reply to agents with cached answers (86% cost reduction)
<b>Question Grouping</b>	Fan-out single answer to multiple waiting agents
<b>Process Hydration</b>	Serialize waiting agents to disk, resume on answer
<b>Resource Locking</b>	Prevent race conditions on shared resources
<b>Cycle Detection</b>	Prevent deadlocks from circular dependencies

### 71.2 Accessing the Dashboard

Navigate to **Blackboard** in the Admin Dashboard sidebar to access:

1. **Overview** - System explanation and architecture benefits
2. **Resolved Facts** - Previously answered questions with reuse counts
3. **Question Groups** - Pending groups waiting for a single answer
4. **Agents** - Active and hydrated agents
5. **Resource Locks** - Currently held locks
6. **Configuration** - System settings

### 71.3 Resolved Facts (Decisions)

The Facts tab shows all resolved questions that can be reused:

Column	Description
<b>Question</b>	The original question asked
<b>Answer</b>	The cached answer
<b>Source</b>	How the answer was obtained (user, memory, default, inferred)
<b>Reused</b>	Number of times this answer has been reused
<b>Status</b>	Valid or Invalid

#### Actions:

- **Invalidate** - Mark an answer as incorrect; affected agents are notified
- **Provide New Answer** - Replace with correct answer during invalidation

### 71.4 Question Grouping

When multiple agents ask similar questions within the grouping window (default: 60 seconds), they are grouped together:

1. First agent's question becomes the "canonical" question
2. Similar questions ( 85% cosine similarity) join the group
3. User answers once
4. Answer is fanned out to all grouped agents

#### Benefits:

- Reduces user interruptions by 60-80%
- Prevents duplicate HITL decisions
- Agents don't wait for individual answers

### 71.5 Process Hydration

Long-running agents are automatically serialized to disk when waiting:

Setting	Default	Description
<b>Auto Hydration</b>	Enabled	Automatically serialize waiting agents
<b>Threshold</b>	300 seconds	Wait time before hydration
<b>S3 Storage</b>	Enabled	Store large states in S3
<b>Compression</b>	gzip	Reduce storage size

#### Agent Lifecycle:

1. Agent starts → **Active**
2. Agent waits for user → **Waiting**
3. Wait exceeds threshold → **Hydrated** (state saved, process killed)
4. User answers → **Restored** (state loaded, execution resumes)

### 71.6 Resource Locking

Prevent race conditions when agents access shared resources:

Lock Type	Description
<b>Read</b>	Multiple readers allowed
<b>Write</b>	Exclusive access, blocks other writers
<b>Exclusive</b>	No other access allowed

**Force Release:** In emergencies, administrators can force-release locks. Use with caution as this may cause data inconsistencies.

## 71.7 Cycle Detection

The system automatically detects circular dependencies:

```
Agent A waits for Agent B
Agent B waits for Agent C
Agent C waits for Agent A ← CYCLE DETECTED
```

### Action on Detection:

1. Cycle is logged with full dependency chain
2. Oldest dependency is broken
3. Affected agent receives timeout error
4. Alert is raised in dashboard

## 71.8 Configuration Options

Setting	Default	Description
similarity_threshold	0.85	Minimum cosine similarity for matching
enable_question_grouping	true	Group similar questions
grouping_window_seconds	60	Wait time for similar questions
enable_answer_reuse	true	Auto-reply with cached answers
answer_ttl_seconds	3600	Answer expiry time
enable_auto_hydration	true	Auto-serialize waiting agents
hydration_threshold_seconds	300	Wait time before hydration
enable_cycle_detection	true	Detect dependency cycles
max_dependency_depth	10	Maximum dependency chain depth

## 71.9 API Endpoints

Base: /api/admin/blackboard

Method	Endpoint	Description
GET	/dashboard	Dashboard statistics
GET	/decisions	List resolved decisions
POST	/decisions/{id}/invalidate	Invalidate a decision
GET	/groups	Pending question groups
POST	/groups/{id}/answer	Answer a group
GET	/agents	Active agents
POST	/agents/{id}/restore	Restore hydrated agent
GET	/locks	Active resource locks
POST	/locks/{id}/release	Force release a lock
GET	/config	Configuration
PUT	/config	Update configuration
POST	/cleanup	Cleanup expired resources
GET	/events	Audit log

## 71.10 Troubleshooting



Issue	Cause	Resolution
Low reuse rate	Threshold too high	Lower <code>similarity_threshold</code> to 0.80
Too many groups	Window too long	Reduce <code>grouping_window_seconds</code>
Agents stuck in hydrated state	Answer never received	Manually restore or let expire
Lock contention	Long-running operations	Increase <code>default_lock_timeout_seconds</code>
Cycle detected frequently	Circular workflows	Review agent dependencies, refactor

### 71.11 Implementation Files

File	Purpose
<code>lambda/shared/services/semantic-blackboard.service.ts</code>	Core service
<code>lambda/shared/services/agent-orchestrator.service.ts</code>	Agent registry, locks
<code>lambda/shared/services/process-hydration.service.ts</code>	State serialization
<code>lambda/admin/blackboard.ts</code>	Admin API
<code>migrations/158_semantic_blackboard_orchestration.sql</code>	Database schema
<code>apps/admin-dashboard/app/(dashboard)/blackboard/page.tsx</code>	Admin UI

## 72. Services Layer & Interface-Based Access Control

The Services Layer ensures all access to RADIANT resources goes through defined interfaces (API, MCP, A2A) with proper authentication and authorization.

### 72.1 Overview

Interface	Purpose	Authentication
<b>API</b>	REST/HTTP endpoints for applications	API Key or JWT
<b>MCP</b>	Model Context Protocol for AI tools	API Key + Capability negotiation
<b>A2A</b>	Agent-to-Agent communication	API Key + mTLS (required by default)

### 72.2 API Keys with Interface Types

Navigate to **API Keys** in the Admin Dashboard to manage keys per interface type.

#### Creating a Key:

1. Click **Create Key**
2. Enter a descriptive name
3. Select **Interface Type**: API, MCP, A2A, or All
4. Configure scopes and expiration
5. For A2A: optionally specify agent ID and type
6. For MCP: optionally restrict allowed tools
7. Click **Create Key**
8. **Copy the key immediately** - it will not be shown again

#### Key Prefixes by Interface:

- `rad_api_` - API-only keys
- `rad_mcp_` - MCP-only keys
- `rad_a2a_` - A2A-only keys
- `rad_all_` - All-interface keys

## 72.3 A2A Agent Management

The A2A Agents tab shows all registered external agents:

Column	Description
<b>Agent</b>	Name and unique ID
<b>Type</b>	Agent category (orchestrator, worker, etc.)
<b>Status</b>	active, suspended, revoked, pending
<b>Operations</b>	Supported A2A operations
<b>Requests</b>	Total request count
<b>Last Heartbeat</b>	Most recent agent heartbeat

### Agent Actions:

- **Suspend:** Temporarily disable agent access
- **Activate:** Re-enable suspended agent
- **Revoke:** Permanently revoke agent access

## 72.4 Interface Policies

Configure access rules per interface in the Policies tab:

Setting	Default	Description
<code>require_authentication</code>	true	Require valid API key
<code>require_mtls</code>	true (A2A)	Require mTLS certificate
<code>global_rate_limit_per_minute</code>	-	Interface-wide rate limit
<code>a2a_require_registration</code>	true	Agents must register first
<code>a2a_max_concurrent_connections</code>	100	Max concurrent A2A connections
<code>mcp_max_tools_per_request</code>	50	Max tools per MCP request

## 72.5 Key Sync Between Admin Apps

Keys created in either Radiant Admin or Think Tank Admin are automatically synced:

1. Key created in Radiant Admin → Queued for Think Tank Admin
2. Sync job processes pending items
3. Key appears in both admin apps

### Sync Status:

- **Pending:** Awaiting sync
- **Synced:** Successfully synchronized
- **Failed:** Sync failed (will retry)

## 72.6 Database Access Restrictions

**CRITICAL:** No external agent can access RADIANT databases directly.

Cedar policies enforce that:

- All database operations require Service principal with `internal=true`
- API keys (api, mcp, a2a, all) are **forbidden** from direct DB access
- External agents must use the defined interfaces

## 72.7 API Endpoints

Base: /api/admin/api-keys

Method	Endpoint	Description
GET	/dashboard	Summary by interface type
GET	/	List all keys
POST	/	Create key with interface type
GET	/:keyId	Get key details
PATCH	/:keyId	Update key
DELETE	/:keyId	Revoke key
POST	/:keyId/restore	Restore revoked key
GET	/agents	List A2A agents
PATCH	/agents/:id/status	Update agent status
GET	/policies	Get interface policies
PUT	/policies/:type	Update policy
GET	/audit	Get audit log
POST	/sync	Process pending syncs

## 72.8 Implementation Files

File	Purpose
migrations/V2026_01_24_001__services_layer_api_keys.sql	Database schema
lambda/admin/api-keys.ts	Admin API handler
lambda/gateway/a2a-worker.ts	A2A protocol worker
config/cedar/interface-access-policies.cedar	Cedar access policies
apps/admin-dashboard/app/(dashboard)/api-keys/page.tsx	Radiant Admin UI
apps/thinktank-admin/app/(dashboard)/api-keys/page.tsx	Think Tank Admin UI

## Section 73: Cortex Graph-RAG Knowledge Engine

**Location:** Admin Dashboard → Memory → Graph-RAG

Enterprise knowledge graph with vector embeddings for intelligent retrieval-augmented generation.

### 73.1 Overview

The Cortex Graph-RAG system provides persistent knowledge storage with semantic relationships:

Component	Technology	Purpose
<b>Entities</b>	PostgreSQL + pgvector	Knowledge nodes with embeddings
<b>Relationships</b>	PostgreSQL	Typed connections between entities
<b>Chunks</b>	PostgreSQL + pgvector	Text segments for RAG retrieval
<b>Vector Search</b>	HNSW Index	Fast approximate nearest neighbor

### 73.2 Dashboard Features

**Overview Page** (/cortex/graph-rag):

- Entity/relationship/chunk counts

- Graph health status (embedding service, vector index)
- Recent activity log
- Top accessed entities

Tabs:

- **Entities:** Browse, search, create, delete knowledge entities
- **Activity:** Recent operations log
- **Configuration:** Feature toggles and retrieval settings

### 73.3 Entity Types

Type	Description	Example
person	Human entities	"John Smith"
organization	Companies, groups	"Acme Corp"
concept	Abstract ideas	"Machine Learning"
event	Occurrences	"Q4 2025 Launch"
location	Places	"San Francisco HQ"
document	Source documents	"User Manual v3"
topic	Subject areas	"Hydraulics"

### 73.4 Relationship Types

Type	Description	Example
is_a	Type hierarchy	"Python is_a Programming Language"
part_of	Composition	"Chapter 1 part_of Manual"
related_to	General association	"AI related_to Machine Learning"
works_for	Employment	"John works_for Acme"
located_in	Geographic	"HQ located_in California"
created_by	Authorship	"Report created_by Jane"
depends_on	Dependency	"Feature depends_on API"

### 73.5 Configuration Options

Setting	Default	Description
enableGraphRag	true	Enable knowledge graph retrieval
enableEntityExtraction	true	Auto-extract entities from content
enableRelationshipInference	true	Auto-infer entity relationships
enableAutoMerge	true	Merge duplicate entities
embeddingModel	text-embedding-3-small	Model for vector embeddings
entityExtractionModel	gpt-4o-mini	Model for entity extraction
defaultMaxResults	10	Max results per query
defaultMaxDepth	3	Max graph traversal depth
minRelevanceScore	0.7	Minimum similarity score
hybridSearchAlpha	0.5	Weight for hybrid search

### 73.6 Vector Search

The system uses HNSW (Hierarchical Navigable Small World) indexing for fast vector similarity search:

```
-- Similarity search function
SELECT * FROM search_cortex_entities(
  p_tenant_id := 'tenant-uuid',
  p_embedding := '[vector]',
  p_limit := 10,
  p_min_similarity := 0.7
);
```

**Hybrid Search** combines:

- Vector similarity (configurable weight)
- Full-text search on name/description
- Graph traversal for relationship context

### 73.7 Graph Traversal

Retrieve entity neighbors using recursive CTE:

```
SELECT * FROM get_entity_neighbors(
  p_tenant_id := 'tenant-uuid',
  p_entity_id := 'entity-uuid',
  p_depth := 2,
  p_relationship_types := ARRAY['related_to', 'part_of']
);
```

### 73.8 Content Ingestion

Ingest content to automatically extract entities and relationships:

POST /api/admin/cortex/ingest

```
{
  "tenantId": "uuid",
  "source": {
    "type": "text",
    "content": "John Smith works at Acme Corp in San Francisco..."
  },
  "options": {
    "extractEntities": true,
    "extractRelationships": true,
    "createChunks": true
  }
}
```

### 73.9 API Endpoints

Base: /api/admin/cortex

Method	Endpoint	Description
GET	/dashboard	Full dashboard data
GET	/config	Get configuration
PUT	/config	Update configuration
GET	/entities	List entities
POST	/entities	Create entity
GET	/entities/:id	Get entity
PUT	/entities/:id	Update entity
DELETE	/entities/:id	Delete entity

Method	Endpoint	Description
GET	/entities/:id/neighbors	Get neighbors
GET	/relationships	List relationships
POST	/relationships	Create relationship
DELETE	/relationships/:id	Delete relationship
GET	/chunks	List chunks
POST	/search	Full-text search
POST	/query	Vector similarity search
POST	/ingest	Ingest content
POST	/merge	Merge entities
GET	/stats	Get statistics

### 73.10 Implementation Files

File	Purpose
packages/shared/src/types/cortex-graph-rag.types.ts	Type definitions
packages/infrastructure/migrations/V2026_01_25_008__cortex_graph_rag.sql	Database schema
packages/infrastructure/lambda/admin/cortex-graph-rag.ts	Admin API
apps/admin-dashboard/app/(dashboard)/cortex/graph-rag/page.tsx	Dashboard UI

### 73.11 Database Tables

Table	Purpose
cortex_config	Per-tenant configuration
cortex_entities	Knowledge entities with embeddings
cortex_relationships	Entity relationships
cortex_chunks	Text chunks for RAG
cortex_activity_log	Activity tracking
cortex_query_log	Query analytics

All tables have RLS policies for multi-tenant isolation using `app.current_tenant_id`.

## Section 75: Complete Admin API Architecture (v5.52.6)

### 75.1 Overview

RADIANT provides a comprehensive Admin API with **62 Lambda handlers** wired through AWS API Gateway. All admin endpoints require Cognito authentication and are protected by admin-level authorization.

### 75.2 Handler Categories

Category	Count	Handlers
<b>Cato Safety</b>	5	cato, cato-genesis, cato-global, cato-governance, cato-pipeline
<b>Memory Systems</b>	4	cortex, cortex-v2, blackboard, empiricism-loop
<b>AI/ML</b>	7	brain, cognition, ego, raws, inference-components, formal-reasoning, ethics-free-reasoning
<b>Security</b>	5	security, security-schedules, api-keys, ethics, self-audit

Category	Count	Handlers
<b>Operations</b>	5	gateway, sovereign-mesh, sovereign-mesh-performance, sovereign-mesh-scaling, hitl-orchestration
<b>Reporting</b>	4	reports, ai-reports, dynamic-reports, metrics
<b>Configuration</b>	7	tenants, invitations, library-registry, checklist-registry, collaboration-settings, system, system-configuration
<b>Infrastructure</b>	6	aws-costs, aws-monitoring, s3-storage, code-quality, infrastructure-tier, logs
<b>Compliance</b>	4	regulatory-standards, council, user-violations, approvals
<b>Models</b>	5	models, lora-adapters, pricing, specialty-rankings, sync-providers
<b>Orchestration</b>	2	orchestration-methods, orchestration-user-templates
<b>Users</b>	2	user-registry, white-label
<b>Time &amp; Translation</b>	3	time-machine, translation, internet-learning
<b>Learning</b>	1	agi-learning

### 75.3 API Route Pattern

All admin handlers are accessible via:

`/api/admin/{handler-name}/*`

Example routes:

- `/api/admin/cato/dashboard` - Cato safety dashboard
- `/api/admin/cortex/memories` - Memory management
- `/api/admin/brain/dreams` - Dream state management
- `/api/admin/tenants/list` - Tenant list

### 75.4 Authentication

All admin endpoints require:

1. **Cognito JWT** - Valid authentication token
2. **Admin Role** - User must have admin privileges
3. **Tenant Context** - Tenant ID for RLS enforcement

### 75.5 Implementation

**CDK Stack:** `packages/infrastructure/lib/stacks/api-stack.ts`

Each handler follows the pattern:

```
const handlerLambda = this.createLambda(
  'HandlerName',
  'admin/handler-name.handler',
  commonEnv, vpc, apiSecurityGroup, lambdaRole
);
const resource = admin.addResource('handler-name');
resource.addProxy({
  defaultIntegration: new apigateway.LambdaIntegration(handlerLambda),
  defaultMethodOptions: {
    authorizer: adminAuthorizer,
    authorizationType: apigateway.AuthorizationType.COGNITO,
  },
});
```

### 75.6 Gap Resolution (v5.52.6)

In v5.52.6, a critical infrastructure audit identified that only ~31 of 62 admin Lambda handlers were wired to API Gateway. The remaining handlers existed but were returning 404 errors. All 62 handlers are now

properly connected.

---

## Section 77: Expert System Adapters (v5.52.21)

### 77.1 Overview

Expert System Adapters (ESA) enable tenant-trainable domain intelligence through automatic LoRA adapter training. Unlike generic AI platforms that treat all tenants the same, ESA allows each tenant to build specialized AI expertise that continuously improves through interaction feedback.

**Key Benefit:** Zero ML expertise required—the system learns automatically from user interactions.

### 77.2 Tri-Layer Adapter Architecture

ESA implements a four-layer adapter stacking system:

$$W_{\text{Final}} = W_{\text{Genesis}} + (\text{scale} \times W_{\text{Cato}}) + (\text{scale} \times W_{\text{User}}) + (\text{scale} \times W_{\text{Domain}})$$

Layer	Name	Purpose	Management
0	Genesis	Base model weights	Frozen, never modified
1	Cato	Global constitution, tenant values	Pinned in memory, never evicted
2	User	Personal preferences, style	LRU eviction when memory constrained
3	Domain	Specialized expertise	Auto-selected by domain detection

### 77.3 Admin Dashboard

**Location:** /models/lora-adapters

The LoRA Adapters page provides:

- **Summary Cards:** Global adapters, User adapters, Invocations (24h), Average latency
- **Tri-Layer Diagram:** Visual representation of adapter stacking
- **Configuration Tab:** Enable/disable adapters, scale settings, auto-selection
- **Adapters Tab:** Registry by layer with activation toggles
- **Warmup Tab:** Manual warmup triggers and history

### 77.4 Configuration Options

Setting	Default	Description
Enable LoRA Adapters	Off	Master toggle for adapter stacking
Use Global Adapter (Cato)	On	Include global constitution adapter
Use User Adapter	On	Include personal preference adapter
Global Scale	1.0	Scaling factor for global adapter (0-2)
User Scale	1.0	Scaling factor for user adapter (0-2)
Auto Selection	Off	Automatically select best adapter
Rollback Enabled	On	Fall back to base model on failure
LRU Eviction	On	Evict least-recently-used adapters
Max Adapters in Memory	50	Maximum loaded adapters
Warmup Interval	15 min	How often to pre-load adapters



## 77.5 Implicit Feedback Learning

ESA automatically captures 11 feedback signals from user behavior:

Signal	Weight	Meaning
Copy Response	+0.80	User copied the output
Thumbs Up	+1.00	Explicit positive feedback
Follow-up Question	+0.30	Partial success, needs more
Long Dwell Time	+0.40	User engaged with response
Share Response	+0.50	User shared the output
Save Response	+0.50	User bookmarked output
Regenerate Request	-0.50	Response wasn't satisfactory
Rephrase Question	-0.50	Original missed the mark
Quick Dismiss	-0.40	User quickly moved on
Abandon Conversation	-0.70	Complete failure
Thumbs Down	-1.00	Explicit negative feedback

## 77.6 Training Pipeline

Training occurs automatically based on configurable thresholds:

Setting	Default	Description
Min Candidates	25	Total candidates before training
Min Positive	15	Minimum positive examples
Min Negative	5	Minimum negative examples
Training Frequency	Weekly	How often to train
Auto Optimal Time	On	Detect best training time

## 77.7 Domain Auto-Selection

When a query arrives, ESA scores available adapters:

$$\text{Score} = (0.3 \times \text{DomainMatch}) + (0.1 \times \text{SubdomainBonus}) + (0.25 \times \text{SatisfactionScore}) \\ + (0.1 \times \text{VolumeScore}) + (0.05 \times \text{ErrorRate}) + (0.2 \times \text{RecencyScore})$$

Adapter selected if Score  $\geq$  0.5

## 77.8 API Endpoints

Base Path: /api/admin/learning

Endpoint	Method	Purpose
/config	GET	Get learning configuration
/config	PUT	Update learning configuration
/domain-adapters	GET	List domain adapters
/domain-adapters/{domain}	GET	Get active adapter for domain
/training/queue	GET	View training queue status
/training/trigger	POST	Manually trigger training
/performance/{adapterId}	GET	Get adapter performance metrics

## 77.9 Implementation Files

File	Purpose
migrations/108_enhanced_learning.sql	Database schema
lambda/shared/services/enhanced-learning.service.ts	Core learning service
lambda/shared/services/lora-inference.service.ts	Tri-layer inference
lambda/shared/services/adapter-management.service.ts	Adapter selection
lambda/admin/enhanced-learning.ts	Admin API handler
apps/admin-dashboard/app/(dashboard)/models/lora-adapters/page.tsx	Admin UI
docs/EXPERT-SYSTEM-ADAPTERS.md	Strategic vision document

## Section 76: PostgreSQL Scaling Infrastructure (v5.52.20)

### 76.1 Overview

**Problem:** When 6 AI models execute in parallel per request, each Lambda opens a database connection. At 100 concurrent requests  $\times$  6 parallel writes = 600 connections—exceeding Aurora's limits and causing transaction conflicts.

**Solution:** OpenAI-inspired PostgreSQL scaling patterns deployed automatically for Tier 2+ installations.

### 76.2 Components

Component	Purpose	Tier
<b>RDS Proxy</b>	Connection pooling, Lambda cold-start optimization	2+
<b>Async Write Queue</b>	SQS-based batch writes for model results	2+
<b>Redis Hot-Path Cache</b>	Read-after-write consistency, rate limiting	2+
<b>Time-Based Partitioning</b>	Monthly partitions for logs/usage tables	All
<b>Materialized Views</b>	Pre-computed dashboard metrics	All
<b>Optimized RLS Policies</b>	Index-friendly tenant isolation	All

### 76.3 RDS Proxy Configuration

Connection limits are tier-based for optimal resource usage:

Tier	Max Connections %	Idle Timeout	Use Case
1	60%	1800s	Development
2	70%	1800s	Starter
3	80%	1200s	Growth
4	85%	900s	Scale
5	90%	600s	Enterprise

### 76.4 Async Write Pattern

Model execution results are written asynchronously to avoid blocking request latency:

Request Flow:

User  $\rightarrow$  Lambda  $\rightarrow$  6 AI Models (parallel)  $\rightarrow$  SQS Queue  $\rightarrow$  Batch Writer  $\rightarrow$  PostgreSQL  
 $\downarrow$   
 Redis Cache (immediate read-after-write)

Queue Configuration:

- Encrypted with tenant KMS key
- 14-day message retention
- 300-second visibility timeout
- Dead letter queue after 3 failures

#### Batch Writer Lambda:

- Processes up to 100 messages per batch
- Tier-based concurrency (5-100 concurrent executions)
- Bulk INSERT for 10-50x efficiency

### 76.5 Redis Caching

Hot-path operations use Redis for immediate consistency:

Operation	Cache TTL	Purpose
Model results	1 hour	Read-after-write consistency
Rate limits	Per window	Tenant/resource throttling
Session state	30 min	Lambda-to-Lambda context

#### Tier-Based Cluster Sizing:

Tier	Node Type	Shards	Replicas
1	cache.t4g.micro	1	0
2	cache.t4g.small	1	1
3	cache.r6g.large	2	1
4	cache.r6g.xlarge	3	2
5	cache.r6g.2xlarge	5	2

### 76.6 Time-Based Partitioning

High-volume tables are partitioned by month:

Table	Partition Key	Retention
model_execution_logs_partitioned	created_at	24 months
usage_records_partitioned	timestamp	24 months

Partitions are:

- Auto-created 3 months ahead
- Archived after 24 months
- Managed via `manage_time_partitions()` function

### 76.7 Materialized Views

Dashboard metrics are pre-computed on schedule:

View	Refresh	Dashboard Use
tenant_daily_usage_summary	15 min	Usage cards
model_performance_summary	1 hour	Model health
tenant_cost_summary	1 hour	Billing

View	Refresh	Dashboard Use
user_activity_summary	1 hour	Engagement
platform_health_stats	5 min	Admin overview
model_popularity_ranking	1 hour	Model selection

## 76.8 Monitoring

Critical metrics to watch:

Metric	Warning	Critical	Action
RDS Proxy connections	< 20%	< 10%	Reduce Lambda concurrency
Aurora CPU	> 70%	> 80%	Add read replicas
SQS queue age	> 30s	> 60s	Increase batch writer concurrency
Query P95 latency	> 300ms	> 500ms	Optimize queries
Redis memory	> 70%	> 80%	Scale cluster

## 76.9 Implementation Files

**CDK Constructs:** | File | Purpose | |-----|-----| | lib/constructs/database-scaling.construct.ts | RDS Proxy CDK construct | | lib/constructs/async-write.construct.ts | SQS + batch writer CDK construct | | lib/constructs/redis-cache.construct.ts | ElastiCache Redis CDK construct | | lib/stacks/data-stack.ts | Integration (tier 2+) |

**Lambda Handlers & Services:** | File | Purpose | |-----|-----| | lambda/scaling/batch-writer.ts | Batch writer Lambda handler with partial failure reporting | | lambda/scaling/model-result-cache.service.ts | Redis cache service for read-after-write | | lambda/scaling/postgresql-scaling.service.ts | Application-level PostgreSQL scaling orchestration |

**Database Migrations (5 total):** | Migration | Purpose | |-----|-----| | V2026\_01\_25\_001\_\_postgresql\_scaling\_rls.sql | Optimized RLS with SELECT wrapper; batch staging; rate limiting | | V2026\_01\_25\_002\_\_postgresql\_scaling\_partitioning.sql | Monthly partitioning; partition management functions | | V2026\_01\_25\_003\_\_postgresql\_scaling\_materialized\_views.sql | 6 materialized views; refresh orchestration | | V2026\_01\_25\_004\_\_postgresql\_scaling\_strategic\_indexes.sql | BRIN/GIN/covering indexes; slow query tracking; index health | | V2026\_01\_25\_005\_\_postgresql\_scaling\_read\_replica\_routing.sql | Read replica routing; session affinity; hot/cold paths |

## 76.10 Strategic Indexing

Index Type	Tables	Purpose
<b>BRIN</b>	model_execution_logs_partitioned, usage_records_partitioned	100x smaller than B-tree for time-series data
<b>Partial</b>	All log tables	Hot-path queries (recent, pending, failed)
<b>Covering</b>	Dashboard queries	Index-only scans eliminate table lookups
<b>GIN</b>	JSONB metadata columns	Efficient containment queries
<b>Expression</b>	Date truncation	Pre-computed daily/hourly grouping

## 76.11 Read Replica Routing

Automatic query routing based on type and consistency requirements:

Query Type	Target	Consistency
Writes	Primary	Strong
Reads (within 5s of write)	Primary	Strong (session affinity)

Query Type	Target	Consistency
Dashboard reads	Any replica	Eventual
Analytics queries	Dedicated replica	Eventual
Materialized view reads	Any replica	Eventual

## 76.12 Slow Query Tracking

Automatic capture of queries exceeding 500ms:

- Query hash for deduplication
- Execution plan capture
- Aggregated statistics in `query_performance_hints`
- Index suggestions via `suggest_indexes()` function

## 76.13 Maintenance Functions

Function	Schedule	Purpose
<code>refresh_priority_materialized_views()</code>	Every 15 min	Dashboard metrics
<code>refresh_all_materialized_views()</code>	Every hour	All materialized views
<code>ensure_future_partitions(3)</code>	Daily	Create next 3 months of partitions
<code>perform_scheduled_maintenance()</code>	Daily	Vacuum, analyze, cleanup
<code>analyze_index_health()</code>	Weekly	Identify unused/inefficient indexes

## 76.14 Admin Dashboard UI

The PostgreSQL Scaling monitoring dashboard is available at:

`/infrastructure/postgresql-scaling`

**Dashboard Tabs:**

Tab	Features
<b>Overview</b>	Connection history, materialized view status, real-time refresh controls
<b>Queues</b>	Batch writer queue status, pending/processing/failed/completed counts, retry failed button
<b>Replicas</b>	Read replica health, lag monitoring, primary/replica status, weights
<b>Partitions</b>	Partition statistics per table, row counts, sizes, ensure future partitions button
<b>Slow Queries</b>	Top slow query patterns, index suggestions, recent slow queries with timing
<b>Maintenance</b>	Manual maintenance triggers, scheduled task overview, maintenance history

**Summary Cards:**

- **Connections:** Active/Max with utilization percentage
- **Queue Status:** Pending count with health indicator
- **Replicas:** Healthy/Total with overall health badge
- **Query Latency:** P95 latency with P50/P99 breakdown

**Admin API Endpoints** (Base: `/api/admin/scaling`):

Endpoint	Method	Purpose
<code>/dashboard</code>	GET	Complete dashboard data
<code>/connections</code>	GET	Connection pool metrics
<code>/queues</code>	GET	Batch writer queue status

Endpoint	Method	Purpose
/queues/retry-failed	POST	Retry failed batch writes
/queues/clear-completed	DELETE	Clear completed writes
/replicas	GET	Read replica health
/partitions	GET	Partition statistics
/partitions/ensure-future	POST	Create future partitions
/slow-queries	GET	Slow query analysis
/indexes	GET	Index health analysis
/indexes/suggestions	GET	Index suggestions
/materialized-views	GET	MV status
/materialized-views/refresh	POST	Trigger MV refresh
/tables	GET	Table statistics
/maintenance/run	POST	Run maintenance
/maintenance/history	GET	Maintenance history
/rate-limits	GET	Rate limiting status

## Version History

Version	Date	Changes
<b>5.52.26</b>	2026-01-25	OAuth 2.0 Provider & Developer Portal (PROMPT-41A); RFC 6749 compliant authorization server
<b>5.52.22</b>	2026-01-25	PostgreSQL Scaling Admin Dashboard; Full visibility into queues, connections, replicas, partitions,
<b>5.52.21</b>	2026-01-25	Expert System Adapters documentation; Strategic vision document; Moat #6D documentation; Tri
<b>5.52.20</b>	2026-01-25	PostgreSQL Scaling Infrastructure; RDS Proxy for connection pooling; Async write pattern with SC
<b>5.52.6</b>	2026-01-24	Complete CDK Wiring Audit; ALL 62 admin Lambda handlers now wired to API Gateway includin
<b>5.52.5</b>	2026-01-24	Services Layer Implementation; API Keys with interface types (API, MCP, A2A); A2A Protocol W
<b>5.52.4</b>	2026-01-24	Semantic Blackboard Admin Dashboard; CDK API route for blackboard Lambda; Complete admin
<b>5.52.0</b>	2026-01-23	Comprehensive UI Audit; Replaced mock data with real API calls; Fixed 17 console.log stubs in Ma
<b>5.46.0</b>	2026-01-23	Cortex Memory System v4.20.0; Three-tier memory architecture (Hot/Warm/Cold); Graph-RAG k
<b>5.42.0</b>	2026-01-22	AI Reports API complete implementation; OpenAPI 3.1 spec ( <a href="#">docs/api/openapi-admin.yaml</a> ); P
<b>5.39.0</b>	2026-01-21	Schema-Adaptive Reports; Dynamic Report Builder; Cato Genesis page; Think Tank Admin naviga
<b>5.38.0</b>	2026-01-21	Sovereign Mesh Performance Optimization; Infrastructure Scaling (100-500K sessions); SQS dispatc
<b>5.37.0</b>	2026-01-21	RAWS v1.1 - Model Selection System; 13 weight profiles; 7 domains with compliance; 8-dimension s
<b>5.34.0</b>	2026-01-20	HITL Orchestration Extensions; Semantic Deduplication with pgvector embeddings; Scout HITL In
<b>5.33.0</b>	2026-01-20	HITL Orchestration Enhancements (PROMPT-37); SAGE-Agent Bayesian VOI; MCP Elicitation S
<b>5.32.0</b>	2026-01-20	Sovereign Mesh Completion; Unit tests for notification and snapshot services; Think Tank Admin in
<b>5.31.0</b>	2026-01-20	The Sovereign Mesh (PROMPT-36); Agent Registry with OODA execution; App Registry (3,000+
<b>5.30.0</b>	2026-01-20	Code Quality & Test Coverage Visibility; Delight service unit tests; JSON safety migration tracking
<b>5.29.0</b>	2026-01-20	Gateway Admin Controls; Persistent statistics with timestamps; Admin dashboard UI; Think Tank
<b>5.28.0</b>	2026-01-20	Multi-Protocol Gateway v3.0; Go Gateway; NATS JetStream; Cedar authorization; Resume tokens
<b>5.20.0</b>	2026-01-18	User Violation Enforcement System; Regulatory compliance tracking; Escalation policies; Appeal w
<b>5.19.0</b>	2026-01-18	White-Label Invisibility (Moat #25); Full branding customization; Custom domains; Response tran
<b>5.7.0</b>	2026-01-17	Deployment Safety (cdk watch DEV-only rule, environment guards, credential setup)
<b>5.6.0</b>	2026-01-12	Genesis Infrastructure (Kaleidos reactor integration, SDS/PDSA compliance, SSF physical-to-digita
<b>5.5.0</b>	2026-01-10	Polymorphic UI (PROMPT-41); ViewRouter component; Terminal/MindMap/DiffEditor views; Ge
<b>5.4.0</b>	2026-01-10	Cognitive Architecture (PROMPT-40); Ghost Memory TTL/semantic key; Economic Governor retr
<b>5.3.0</b>	2026-01-10	Semantic Blackboard; Multi-Agent Orchestration; MCP Primary Interface; Process Hydration; Cyc
<b>5.0.2</b>	2026-01-08	The Grimoire (procedural memory); Economic Governor (cost optimization); Self-optimizing archite
4.20.3	2026-01-08	Mission Control GA; MCP Hybrid Interface; Domain Risk Policies
4.20.2	2026-01-07	Fixed RLS tenant bleed
4.20.0	2026-01-06	Initial Mission Control release

---

*Version 5.52.26 / January 2026 Cross-AI Validated: Claude Opus 4.5 / Google Gemini System Evolution:  
OAuth 2.0 Provider & Developer Portal Status: GO FOR LAUNCH*