

# Contents

<b>RADIANT &amp; Think Tank Executive Summary</b>	<b>6</b>
What is RADIANT?	6
What is Think Tank?	7
Key Differentiators	7
1. AGI-Driven Model Selection	7
2. 49 Proven Orchestration Patterns	8
3. Enterprise-Grade Security	8
Platform Components	8
By the Numbers	8
Use Cases	9
Enterprise AI Gateway	9
Complex Problem Solving (Think Tank)	9
Quality-Critical Applications	9
Cost Optimization	9
Competitive Advantages	9
Technology Stack	10
Deployment Model	10
Pricing Model	10
Roadmap Highlights	11
Summary	11
<b>RADIANT Platform Architecture</b>	<b>11</b>
Table of Contents	12
1. Platform Overview	12
1.1 What is RADIANT?	12
1.2 Core Value Proposition	12
1.3 Platform Statistics	12
2. System Architecture	13
2.1 High-Level Architecture	13
2.2 Three-Component Structure	14
3. Component Deep Dive	15
3.1 Model Router Service	15
3.2 Service Architecture	16
4. Data Architecture	17
4.1 Database Schema Overview	17
4.2 Multi-Tenant Data Isolation	19
5. Security Architecture	20
5.1 Security Layers	20
6. Deployment Architecture	21
6.1 AWS Infrastructure	21
6.2 CDK Stack Dependencies	22
7. Integration Points	23
7.1 External API Integrations	23
<b>Think Tank Platform Architecture</b>	<b>24</b>
Table of Contents	24

1. Platform Overview	24
1.1 What is Think Tank?	24
1.2 Think Tank vs Traditional Chat	24
1.3 Key Capabilities	25
2. Core Architecture	25
2.1 System Components	25
2.2 Think Tank Engine	27
3. Problem Solving Pipeline	28
3.1 Pipeline Stages	28
3.2 Step Recording	30
4. Session Management	31
4.1 Session Lifecycle	31
4.2 Session Data Model	32
5. Collaboration Features	33
5.1 Real-Time Collaboration	33
6. Domain Modes	34
6.1 Specialized Reasoning Modes	34
7. Quality & Confidence	36
7.1 Confidence Scoring System	36
8. User Interface	38
8.1 Think Tank UI Layout	38
<b>AGI &amp; Workflow Orchestration</b>	<b>39</b>
1. Overview	39
Why 50-300% Improvement?	39
Key Capabilities	40
Improvement by Use Case	40
2. The 49 Orchestration Patterns	41
Pattern Categories	41
3. AGI Dynamic Model Selection	42
How It Works	42
Domain Detection Keywords	43
4. Model Execution Modes	43
5. Parallel Execution	44
Execution Modes	44
Synthesis Strategies	44
6. Visual Workflow Editor	44
Editor Features	44
Step Configuration	44
7. API Usage	45
Execute Workflow	45
8. Benefits	45
<b>Simultaneous Prompt Execution</b>	<b>46</b>
Overview	46
RADIANT Parallel Execution	46
Configuration	46
Execution Modes	46

Implementation . . . . .	46
Use Cases . . . . .	48
Think Tank Concurrent Sessions . . . . .	49
Session-Level Parallelism . . . . .	49
Implementation . . . . .	49
Session Configuration . . . . .	50
Database Schema Support . . . . .	51
Performance Benefits . . . . .	51
Throughput Improvement . . . . .	51
Quality Improvement from Consensus . . . . .	51
API Usage . . . . .	52
REST API . . . . .	52
Response . . . . .	52
SDK Usage . . . . .	52
Cost Considerations . . . . .	53
Feature Categories . . . . .	53
1. AI Model Management . . . . .	54
1.1 Model Router Service . . . . .	54
1.2 Model Metadata Service . . . . .	54
1.3 Supported Models (106+) . . . . .	54
2. Orchestration & Workflows . . . . .	55
2.1 Orchestration Patterns (49) . . . . .	55
2.2 AGI Dynamic Model Selection . . . . .	55
2.3 Model Execution Modes (9) . . . . .	55
2.4 Parallel Execution . . . . .	56
2.5 Visual Workflow Editor . . . . .	56
3. Think Tank Platform . . . . .	56
3.1 Problem Solving Engine . . . . .	56
3.2 Session Management . . . . .	57
3.3 Domain Modes (8) . . . . .	57
3.4 Collaboration . . . . .	57
4. Billing & Cost Management . . . . .	57
4.1 Credit System . . . . .	57
4.2 Subscriptions . . . . .	58
4.3 Cost Management . . . . .	58
5. Multi-Tenant Platform . . . . .	58
5.1 Tenant Management . . . . .	58
5.2 User Management . . . . .	58
5.3 API Key Management . . . . .	59
6. Security & Compliance . . . . .	59
6.1 Data Security . . . . .	59
6.2 Authentication . . . . .	59
6.3 Compliance . . . . .	59
7. Analytics & Monitoring . . . . .	59
7.1 Usage Analytics . . . . .	59
7.2 Model Performance . . . . .	60
7.3 Business Intelligence . . . . .	60
8. Developer Tools . . . . .	60

8.1 SDK . . . . .	60
8.2 Webhooks . . . . .	60
8.3 Integrations . . . . .	61
9. Admin Dashboard . . . . .	61
9.1 Dashboard Pages . . . . .	61
9.2 UI Features . . . . .	61
10. Swift Deployer App . . . . .	61
10.1 Deployment Features . . . . .	61
10.2 QA & Testing . . . . .	62
10.3 AI Assistant . . . . .	62
10.4 Local Storage . . . . .	62
<b>RADIANT Services Reference</b>	<b>62</b>
Complete Lambda Services Inventory (62 Services) . . . . .	62
Core Infrastructure Services . . . . .	62
AI Model Services . . . . .	63
Orchestration Services . . . . .	66
Billing Services . . . . .	69
Cognitive Services . . . . .	69
Memory Services . . . . .	70
AGI Services . . . . .	71
Collaboration Services . . . . .	71
Additional Services (30-62) . . . . .	72
<b>RADIANT Database Schema Reference</b>	<b>73</b>
Complete Migration Inventory (40 Migrations) . . . . .	73
Migration Index . . . . .	73
Core Tables (Migration 001) . . . . .	74
tenants . . . . .	74
users . . . . .	75
administrators . . . . .	76
approval_requests . . . . .	76
Think Tank Tables (Migration 016) . . . . .	77
thinktank_sessions . . . . .	77
thinktank_steps . . . . .	77
thinktank_tools . . . . .	78
Orchestration Tables (Migration 024) . . . . .	79
workflow_definitions . . . . .	79
workflow_tasks . . . . .	80
workflow_executions . . . . .	81
task_executions . . . . .	81
Billing Tables (Migrations 033-035) . . . . .	82
credit_balances . . . . .	82
credit_transactions . . . . .	83
subscription_tiers . . . . .	83
subscriptions . . . . .	84
Row-Level Security (Migration 002) . . . . .	84
RLS Pattern . . . . .	84

Setting Tenant Context . . . . .	84
Tables with RLS Enabled: . . . . .	84
Common Patterns . . . . .	85
Updated At Trigger . . . . .	85
Soft Delete Pattern . . . . .	85
Audit Columns . . . . .	85
<b>Swift Deployer App Reference</b> . . . . .	<b>86</b>
App Architecture . . . . .	86
Overview . . . . .	86
File Structure (36 Files) . . . . .	86
Models . . . . .	87
Configuration.swift . . . . .	87
Credentials.swift . . . . .	88
Deployment.swift . . . . .	88
Services . . . . .	89
CDKService.swift . . . . .	89
DeploymentService.swift . . . . .	90
AIAssistantService.swift . . . . .	92
LocalStorageManager.swift . . . . .	93
HealthCheckService.swift . . . . .	95
Components . . . . .	96
MacOSComponents.swift . . . . .	96
AppCommands.swift . . . . .	97
UI Patterns (10 macOS Patterns) . . . . .	99
Dashboard Architecture . . . . .	99
Technology Stack . . . . .	99
Page Structure . . . . .	99
Complete Page Inventory (43 Pages) . . . . .	100
Core Administration . . . . .	100
AI & Models . . . . .	100
Orchestration . . . . .	100
Think Tank . . . . .	100
Billing & Cost . . . . .	101
Analytics & Monitoring . . . . .	101
AGI & Learning . . . . .	101
Collaboration & Features . . . . .	101
Key Pages Detail . . . . .	102
Overview Dashboard (/) . . . . .	102
Models Page (/models) . . . . .	102
Orchestration Patterns Page (/orchestration-patterns) . . . . .	102
Visual Workflow Editor (/orchestration-patterns/editor) . . . . .	103
Billing Page (/billing) . . . . .	103
Analytics Page (/analytics) . . . . .	104
Security Page (/security) . . . . .	104
Think Tank Page (/thinktank) . . . . .	104
Component Library . . . . .	105
Shared Components . . . . .	105

API Integration . . . . .	105
API Client ( <code>lib/api.ts</code> ) . . . . .	105
Authentication ( <code>lib/auth.ts</code> ) . . . . .	106
<b>Compliance &amp; Security Standards</b>	<b>107</b>
Overview . . . . .	107
Required Provider API Keys . . . . .	107
Deployment Flow . . . . .	107
Why These Providers Are Required . . . . .	107
SOC 2 Type II Compliance . . . . .	107
Trust Service Criteria . . . . .	107
Key Controls . . . . .	108
HIPAA Compliance . . . . .	108
Protected Health Information (PHI) Handling . . . . .	108
HIPAA Mode Features . . . . .	108
PHI Sanitization . . . . .	109
GDPR Compliance . . . . .	109
Data Subject Rights . . . . .	109
Data Processing . . . . .	110
Consent Management . . . . .	110
Data Retention . . . . .	111
ISO 27001 Compliance . . . . .	111
Information Security Management System (ISMS) . . . . .	111
Annex A Controls . . . . .	111
Risk Assessment Matrix . . . . .	112
Security Architecture . . . . .	113
Defense in Depth . . . . .	113
Two-Person Approval . . . . .	113
Audit Logging . . . . .	114
Log Structure . . . . .	114
Logged Actions . . . . .	115
Log Retention . . . . .	115

## RADIANT & Think Tank Executive Summary

### Enterprise AI Platform Overview

Version 4.18.0 | December 2024

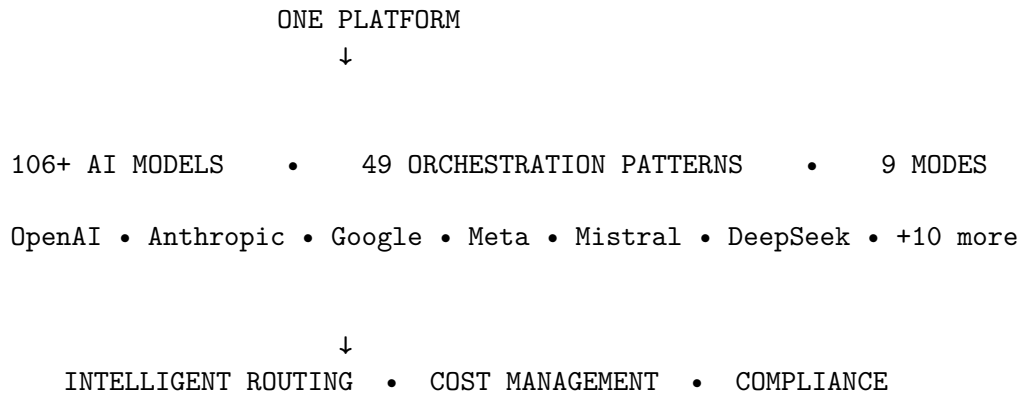
---

*For executives, investors, and decision-makers*

---

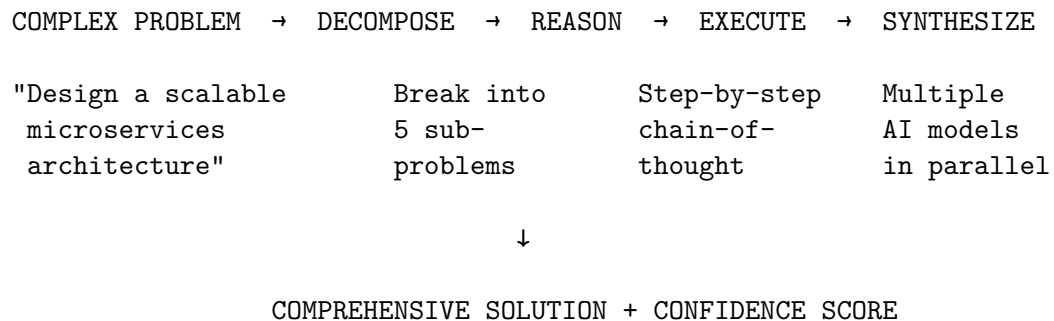
### What is RADIANT?

**RADIANT** is an enterprise-grade, multi-tenant AI platform that provides organizations with unified access to 106+ AI models through a single API, with intelligent orchestration that coordinates multiple AI systems to deliver superior results.



## What is Think Tank?

**Think Tank** is RADIANT's advanced problem-solving platform that decomposes complex problems into manageable steps, applies multi-AI reasoning, and synthesizes comprehensive solutions with confidence scoring.



## Key Differentiators

### 1. AGI-Driven Model Selection

Unlike platforms that use a single AI model, RADIANT's AGI layer **automatically selects the optimal combination of models** based on task analysis:

What We Analyze	What We Select
Problem domain (coding, legal, medical...)	Best models for that domain
Task complexity	Number of models (2-5)

What We Analyze	What We Select
Reasoning requirements	Execution mode (thinking, fast, precise...)
Quality vs speed priority	Parallel execution strategy

**Result:** 50-300% better outcomes than single-model approaches.

## 2. 49 Proven Orchestration Patterns

Research-backed workflows including: - **AI Debate** - Two AIs argue, judge decides - **Self-Refine** - Generate → Critique → Improve - **Chain-of-Verification** - Fact-check every claim - **Tree of Thoughts** - Explore multiple solution paths

## 3. Enterprise-Grade Security

Capability	Description
<b>Multi-Tenant Isolation</b>	PostgreSQL Row-Level Security
<b>Compliance</b>	SOC2, HIPAA-ready
<b>Encryption</b>	At-rest (AES-256) and in-transit (TLS 1.3)
<b>Audit Logging</b>	Complete activity trail

## Platform Components

### RADIANT PLATFORM

SWIFT DEPLOYER (macOS)	AWS INFRASTRUCTURE	ADMIN DASHBOARD (Next.js)
One-click deployment & management	<ul style="list-style-type: none"> <li>• 14 CDK Stacks</li> <li>• Lambda</li> <li>• Aurora PG</li> <li>• API Gateway</li> <li>• S3, Redis</li> </ul>	Manage: <ul style="list-style-type: none"> <li>• Tenants</li> <li>• Users</li> <li>• Billing</li> <li>• Analytics</li> </ul>

## By the Numbers



---

Metric	Value
<b>AI Models</b>	106+ (50 external + 56 self-hosted)
<b>AI Providers</b>	15+ integrated
<b>Orchestration Patterns</b>	49 documented workflows
<b>Execution Modes</b>	9 specialized modes
<b>Database Migrations</b>	66+ schema versions
<b>CDK Stacks</b>	14 infrastructure components

---



---

## Use Cases

### Enterprise AI Gateway

- Unified access to all major AI providers
- Centralized cost management and budgeting
- Consistent API regardless of backend model
- Automatic failover for reliability

### Complex Problem Solving (Think Tank)

- Multi-step technical analysis
- Research synthesis with citations
- Architecture design with artifacts
- Decision support with confidence scores

### Quality-Critical Applications

- Legal document analysis (precise mode)
- Medical information processing (HIPAA compliant)
- Financial analysis (multi-model verification)
- Code generation (AI debate + critique)

### Cost Optimization

- Intelligent model routing (use cheaper models when appropriate)
  - Budget alerts and limits
  - Usage analytics by team/project
  - Model performance vs cost analysis
- 

## Competitive Advantages

---

vs. Single-Model APIs	vs. Other Platforms
Multi-model orchestration	49 research-backed patterns
Built-in verification	AGI-driven model selection
Higher accuracy	9 execution modes

vs. Single-Model APIs	vs. Other Platforms
Reduced bias Confidence scoring	Visual workflow editor Think Tank problem solving

## Technology Stack

Layer	Technology
<b>Frontend</b>	Next.js 14, TypeScript, Tailwind CSS, shadcn/ui
<b>Backend</b>	AWS Lambda (Node.js 20), API Gateway
<b>Database</b>	Aurora PostgreSQL (Serverless), DynamoDB, Redis
<b>Infrastructure</b>	AWS CDK (TypeScript), 14 stacks
<b>Desktop</b>	SwiftUI (macOS 13.0+, Swift 5.9+)
<b>Security</b>	Cognito, KMS, WAF, Row-Level Security

## Deployment Model

RADIANT deploys to **your AWS account**:

YOUR AWS ACCOUNT

RADIANT INFRASTRUCTURE

- Your data stays in your account
- Your compliance requirements met
- Your region/residency requirements
- Full control over infrastructure

Deployed via Swift Deployer (macOS app) or CLI

## Pricing Model

Tier	Target	Includes
<b>Free</b>	Developers	10K tokens/month, 3 models

Tier	Target	Includes
<b>Pro</b>	Teams	1M tokens/month, all models, orchestration
<b>Enterprise</b>	Organizations	Unlimited, SLA, custom patterns, HIPAA

All tiers include: - Full API access - Admin dashboard - Basic analytics - Email support

## Roadmap Highlights

Timeframe	Features
<b>Q1 2026</b>	Mobile SDK, more self-hosted models
<b>Q2 2026</b>	Fine-tuning pipeline, custom model hosting
<b>Q3 2026</b>	Multi-region deployment, advanced compliance
<b>Q4 2026</b>	Marketplace for custom patterns

## Summary

**RADIANT + Think Tank** delivers:

1. **Unified AI Access** - One API for 106+ models across 15+ providers
2. **Intelligent Orchestration** - AGI selects optimal models and modes
3. **Superior Results** - 49 patterns achieve 50-300% better outcomes
4. **Enterprise Security** - Multi-tenant, SOC2, HIPAA-ready
5. **Cost Control** - Budgets, analytics, intelligent routing
6. **Problem Solving** - Think Tank for complex multi-step reasoning

## RADIANT v4.18.0 + Think Tank v3.2.0

*The enterprise platform for intelligent AI orchestration*

**Contact:** [info@radiant.ai](mailto:info@radiant.ai) | **Documentation:** [docs.radiant.ai](https://docs.radiant.ai)

## RADIANT Platform Architecture

### Enterprise Multi-Tenant AI Platform

Version 4.18.0 | December 2024

*A comprehensive technical architecture document for the RADIANT AI orchestration platform*

Table of Contents

- 1. Platform Overview
- 2. System Architecture
- 3. Component Deep Dive
- 4. Data Architecture
- 5. Security Architecture
- 6. Deployment Architecture
- 7. Integration Points

1. Platform Overview

1.1 What is RADIANT?

**RADIANT** (Real-time AI Distribution, Integration, and Automation Network for Tenants) is an enterprise-grade, multi-tenant SaaS platform that provides unified access to 106+ AI models across multiple providers, with intelligent orchestration, cost management, and comprehensive analytics.

1.2 Core Value Proposition

RADIANT VALUE PROPOSITION			
UNIFIED ACCESS	INTELLIGENT ORCHESTRATION	COST MANAGEMENT	ENTERPRISE SECURITY
106+ AI Models One API	49 Multi-AI Patterns AGI Router	Credits, Budgets, Analytics	SOC2/HIPAA Compliant Multi-Tenant

1.3 Platform Statistics

Metric	Value
AI Models Supported	106+ (50 external + 56 self-hosted)
AI Providers Integrated	15+ (OpenAI, Anthropic, Google, Meta, etc.)
Orchestration Patterns	49 documented patterns
Model Execution Modes	9 (thinking, research, fast, creative, etc.)
Database Migrations	66+ schema migrations
CDK Stacks	14 infrastructure stacks

## 2. System Architecture

### 2.1 High-Level Architecture

#### RADIANT PLATFORM ARCHITECTURE

##### CLIENT LAYER

Swift Deployer (macOS)	Admin Dashboard (Next.js)	Think Tank Consumer App	SDK & API Clients
------------------------------	---------------------------------	-------------------------------	-------------------------

##### API GATEWAY LAYER

	Amazon API Gateway	
• REST APIs	• WebSocket APIs	• Rate Limiting
• JWT Auth	• API Keys	• Usage Plans

##### COMPUTE LAYER (Lambda)

Model Router	Orchestration Engine	Billing Service	Admin Service
Think Tank	Analytics Service	Learning Service	Webhook Service

##### DATA LAYER

Aurora PostgreSQL	DynamoDB (Sessions)	S3 (Storage)	Redis (Cache)
----------------------	------------------------	-----------------	------------------

(RLS)

#### EXTERNAL AI PROVIDERS

OpenAI	Anthropic	Google	Meta	Mistral	+10 more
GPT-4o	Claude	Gemini	Llama		
o1	3.5	2.0	3.1		

## 2.2 Three-Component Structure

RADIANT consists of three primary deployment components:

#### THREE COMPONENTS OF RADIANT

packages/infrastructure/lambda/shared/services/

#### CORE SERVICES

model-router.service.ts	Route requests to AI providers
model-metadata.service.ts	Live model data & capabilities
orchestration-patterns.service	49 multi-AI workflow patterns
superior-orchestration.service	Guaranteed superior responses
learning.service.ts	ML feedback & continuous learning

#### BILLING SERVICES

billing.service.ts	Credit & subscription management
cost-management.service.ts	Budget alerts & cost tracking
usage-analytics.service.ts	Usage metrics & reporting

## PLATFORM SERVICES

<code>tenant.service.ts</code>	Multi-tenant management
<code>auth.service.ts</code>	Authentication & authorization
<code>api-key.service.ts</code>	API key lifecycle
<code>webhook.service.ts</code>	Event notifications
<code>storage.service.ts</code>	File & artifact storage

## THINK TANK SERVICES

<code>thinktank-engine.ts</code>	Multi-step problem solving
<code>thinktank-sessions.ts</code>	Conversation management
<code>collaboration.service.ts</code>	Real-time collaboration

---

### 3. Component Deep Dive

#### 3.1 Model Router Service

The intelligent core that routes AI requests to optimal providers:

#### MODEL ROUTER SERVICE

Incoming  
Request

- Request Validation
- API Key Check
  - Rate Limiting
  - Tenant Verification

Model Selection	Budget Check	Fallback Logic
<ul style="list-style-type: none"> <li>• Metadata</li> <li>• Preferences</li> <li>• Capabilities</li> </ul>	<ul style="list-style-type: none"> <li>• Credits</li> <li>• Limits</li> <li>• Cost Est.</li> </ul>	<ul style="list-style-type: none"> <li>• Primary</li> <li>• Secondary</li> <li>• Tertiary</li> </ul>

#### Provider Adapter Layer

OpenAI Adapter	Anthropic Adapter	...
-------------------	----------------------	-----

#### Response Processing

- Token Counting
- Cost Calculation
- Usage Recording
- Analytics Event

## 3.2 Service Architecture

### LAMBDA SERVICES ARCHITECTURE

packages/infrastructure/lambda/shared/services/

#### CORE SERVICES

model-router.service.ts	Route requests to AI providers
model-metadata.service.ts	Live model data & capabilities
orchestration-patterns.service	49 multi-AI workflow patterns



superior-orchestration.service	Guaranteed superior responses
learning.service.ts	ML feedback & continuous learning

#### BILLING SERVICES

billing.service.ts	Credit & subscription management
cost-management.service.ts	Budget alerts & cost tracking
usage-analytics.service.ts	Usage metrics & reporting

#### PLATFORM SERVICES

tenant.service.ts	Multi-tenant management
auth.service.ts	Authentication & authorization
api-key.service.ts	API key lifecycle
webhook.service.ts	Event notifications
storage.service.ts	File & artifact storage

#### THINK TANK SERVICES

thinktank-engine.ts	Multi-step problem solving
thinktank-sessions.ts	Conversation management
collaboration.service.ts	Real-time collaboration

---

## 4. Data Architecture

### 4.1 Database Schema Overview

#### AURORA POSTGRESQL SCHEMA

66+ Migrations in packages/infrastructure/migrations/

## CORE ENTITIES

tenants	Multi-tenant organizations
users	User accounts with roles
api_keys	API authentication keys
model_configurations	Per-tenant model settings
model_metadata	AI model capabilities & pricing

## BILLING & CREDITS

credit_accounts	Tenant credit balances
credit_transactions	Credit usage history
subscriptions	Plan subscriptions
invoices	Billing invoices
budgets	Spending limits & alerts

## ORCHESTRATION

orchestration_methods	Reusable AI method definitions
orchestration_workflows	49 workflow patterns
workflow_method_bindings	Steps linking workflows to methods
orchestration_executions	Execution history & results

## THINK TANK

thinktank_sessions	Problem-solving sessions
thinktank_conversations	Conversation threads
thinktank_messages	Individual messages
thinktank_steps	Reasoning steps
thinktank_artifacts	Generated outputs

## ANALYTICS & LEARNING

usage_events	API usage events
analytics_aggregates	Pre-computed metrics
learning_interactions	ML training data
model_performance	Model quality tracking

## SECURITY

Row-Level Security (RLS) on all tenant tables  
SET app.current\_tenant\_id for automatic filtering  
Audit logging on sensitive operations

### 4.2 Multi-Tenant Data Isolation

#### ROW-LEVEL SECURITY (RLS) MODEL

Request from Tenant A

Request from Tenant B

JWT Token  
tenant\_id=A

JWT Token  
tenant\_id=B

#### Database Connection

```
SET app.current_tenant_id = 'tenant_id_from_jwt';
```

#### RLS Policy Applied

```
CREATE POLICY tenant_isolation ON table_name  
USING (tenant_id = current_setting('app.current_tenant_id'));
```

Result: Each tenant ONLY sees their own data

Tenant A sees:

Tenant B sees:

Only Tenant A's

- Users
- API Keys
- Usage Data

Only Tenant B's

- Users
- API Keys
- Usage Data

---

## 5. Security Architecture

### 5.1 Security Layers

#### SECURITY ARCHITECTURE

##### LAYER 1: NETWORK SECURITY

- VPC with private subnets for database
- WAF rules for API Gateway
- CloudFront for DDoS protection
- TLS 1.3 for all connections

##### LAYER 2: AUTHENTICATION

- Cognito User Pools for user authentication
- JWT tokens with tenant claims
- API Keys with scoped permissions
- MFA support for admin users

##### LAYER 3: AUTHORIZATION

- Role-based access control (RBAC)
- Permission sets per tenant
- Resource-level policies
- API endpoint authorization

##### LAYER 4: DATA SECURITY

- Row-Level Security (RLS) in PostgreSQL
- Encryption at rest (AES-256)
- Encryption in transit (TLS)
- KMS for key management

- PHI sanitization for HIPAA compliance

#### LAYER 5: AUDIT & COMPLIANCE

- CloudTrail for API logging
- Audit tables for data changes
- Compliance reporting dashboard
- SOC2 Type II controls
- HIPAA compliance mode

---

## 6. Deployment Architecture

### 6.1 AWS Infrastructure

#### AWS DEPLOYMENT ARCHITECTURE

REGION: us-east-1

CloudFront CDN

S3 Static Assets

API Gateway  
(REST + WS)

Lambda Functions  
(Node.js 20)

Cognito  
User Pools

Aurora  
PostgreSQL  
(Serverless)

S3  
Storage

Secrets

ElastiCache

SQS

Manager (Redis) Queues

## 6.2 CDK Stack Dependencies

### CDK STACK DEPENDENCIES

NetworkStack  
(VPC, Subnets)

AuthStack DatabaseStack StorageStack  
(Cognito) (Aurora) (S3)

AISStack  
(Model Router)

APIStack BillingStack ThinkTankStack

Additional stacks: AnalyticsStack, WebhookStack, ComplianceStack,  
MonitoringStack, CDNStack, NotificationStack

7. Integration Points

7.1 External API Integrations

EXTERNAL INTEGRATIONS

AI PROVIDERS (15+)

OpenAI GPT-4o, o1	Anthropic Claude 3.5	Google Gemini 2.0	Meta Llama 3.1	Mistral Large
Cohere	AI21	Perplexity Sonar	DeepSeek R1, Chat	xAI Grok

PAYMENT PROVIDERS

Stripe	Credit card processing, subscriptions, invoicing
--------	--

MONITORING & OBSERVABILITY

CloudWatch (Logs)	X-Ray (Traces)	Sentry (Errors)	Custom Analytics Dashboard
----------------------	-------------------	--------------------	-------------------------------

NOTIFICATIONS

SES (Email)	SNS (Push)	Webhooks (Custom)	Slack/Teams Integrations
----------------	---------------	----------------------	--------------------------

---

**RADIANT Platform Architecture v4.18.0**

*Building the future of enterprise AI*

---

# Think Tank Platform Architecture

## Advanced Multi-Step AI Problem Solving

Version 3.2.0 | December 2024

*A comprehensive technical architecture document for the Think Tank AI reasoning platform*

### Table of Contents

- 1. Platform Overview
- 2. Core Architecture
- 3. Problem Solving Pipeline
- 4. Session Management
- 5. Collaboration Features
- 6. Domain Modes
- 7. Quality & Confidence
- 8. User Interface

## 1. Platform Overview

### 1.1 What is Think Tank?

**Think Tank** is an advanced AI reasoning platform that decomposes complex problems into manageable sub-problems, applies multi-step reasoning, and synthesizes comprehensive solutions using orchestrated AI models.

Unlike simple chat interfaces, Think Tank: - **Decomposes** complex problems into sub-tasks - **Reasons** through each component step-by-step - **Executes** specialized AI calls for each step - **Synthesizes** results into coherent solutions - **Tracks** confidence and quality throughout

### 1.2 Think Tank vs Traditional Chat

TRADITIONAL CHAT vs THINK TANK			
TRADITIONAL CHAT		THINK TANK	
User	AI Response	User	Problem Analysis
Single prompt, single response			
No decomposition			
No reasoning steps		Decompose	
No confidence tracking		into parts	



No iterative refinement

Part 1      Part 2      Part 3  
Reason      Reason      Reason

Execute      Execute  
+ Verify      + Verify

Synthesize  
Solution  
(confidence)

1.3 Key Capabilities

Capability	Description
Problem Decomposition	Breaks complex questions into manageable sub-problems
Multi-Step Reasoning	Chain-of-thought with recorded steps
Domain Specialization	8+ specialized reasoning modes
Confidence Tracking	Quality scores for every step
Artifact Generation	Code, documents, diagrams as outputs
Real-time Collaboration	Multiple users solving together
Session Persistence	Resume any session later
Cost Transparency	Token and cost tracking per step

2. Core Architecture

2.1 System Components

THINK TANK ARCHITECTURE

## CONSUMER INTERFACE LAYER

Web Client (Next.js/React)	Mobile Client (React Native)	API Client (SDK)
-------------------------------	---------------------------------	---------------------

## THINK TANK ENGINE

Session Manager	Problem Decomposer	Step Executor
Reasoning Engine	Solution Synthesizer	Confidence Scorer

## ORCHESTRATION LAYER

### OrchestrationPatternsService

- 49 workflow patterns
- Parallel execution
- AGI model selection
- Mode-aware invocation

### ModelRouterService

- 106+ AI models
- Live metadata
- Intelligent routing
- Fallback handling

## DATA LAYER

Sessions (Aurora)	Conversations (Aurora)	Messages (Aurora)	Artifacts (S3)
----------------------	---------------------------	----------------------	-------------------

## 2.2 Think Tank Engine

The core engine that powers intelligent problem solving:

### THINK TANK ENGINE DETAIL

```
class ThinkTankEngine {  
  
    async solve(problem: ThinkTankProblem): Promise<ThinkTankResult>  
  
        1. CREATE SESSION  
            • Initialize session with problem context  
            • Set domain mode and preferences  
            • Record start time and user info  
  
        2. DECOMPOSE PROBLEM  
            • AI analyzes problem structure  
            • Identifies sub-problems and dependencies  
            • Creates execution plan  
  
        3. FOR EACH SUB-PROBLEM:  
  
            a. REASON  
                • Chain-of-thought analysis  
                • Record reasoning steps  
  
            b. EXECUTE  
                • Call appropriate AI model(s)  
                • May use parallel execution  
                • Track tokens and cost  
  
            c. RECORD STEP  
                • Save step result with confidence  
                • Update session state  
  
        4. SYNTHESIZE SOLUTION  
            • Combine all step results  
            • Generate final answer with reasoning  
            • Calculate overall confidence
```

```

5. RETURN RESULT
  • Solution with confidence score
  • All recorded steps
  • Total cost and token usage
}

```

---

### 3. Problem Solving Pipeline

#### 3.1 Pipeline Stages

##### PROBLEM SOLVING PIPELINE

###### USER INPUT

"Design a scalable microservices architecture for an e-commerce platform that handles 10M daily users with real-time inventory"

###### STAGE 1: PROBLEM ANALYSIS

- Identify problem type: System Design
- Detect domain: Engineering/Architecture
- Assess complexity: High
- Select domain mode: Engineering Mode
- Choose orchestration pattern: Decomposed Prompting

###### STAGE 2: DECOMPOSITION

Sub-Problem 1: Requirements Analysis  
 Sub-Problem 2: Service Identification  
 Sub-Problem 3: Data Architecture  
 Sub-Problem 4: Communication Patterns  
 Sub-Problem 5: Scalability Design  
 Sub-Problem 6: Infrastructure

Dependencies: [1] → [2,3] → [4] → [5] → [6]

### STAGE 3: STEP-BY-STEP EXECUTION

#### Step 1: Requirements

Model: Claude 3.5 (thinking mode)

Tokens: 2,450 Cost: \$0.024 Confidence: 0.92

Output: Detailed requirements document

#### Step 2: Service Identification

Model: GPT-4o + Claude (parallel, merge synthesis)

Tokens: 3,200 Cost: \$0.041 Confidence: 0.89

Output: 12 microservices identified with boundaries

[Steps 3-6 continue...]

### STAGE 4: SYNTHESIS

- Combine all step outputs
- Generate comprehensive solution document
- Include architecture diagram (artifact)
- Validate consistency across steps
- Calculate final confidence: 0.88

### FINAL OUTPUT

- Complete microservices architecture document
- Service interaction diagrams
- Database schema recommendations
- Infrastructure as code templates
- Scaling strategies and benchmarks

Total: 12,400 tokens \$0.18 6 steps 45 seconds

### 3.2 Step Recording

Every reasoning step is recorded with comprehensive metadata:

#### STEP RECORD STRUCTURE

```
interface ThinkTankStep {
    stepId: string;           // Unique step identifier
    sessionId: string;        // Parent session
    stepOrder: number;        // Execution order
    stepType: StepType;       // decompose | reason | execute | ..
    title: string;            // Human-readable step name
    description: string;       // What this step does

    // Execution Details
    input: {
        prompt: string;       // Input to AI
        context: Record<string, any>; // Previous step outputs
        parameters: Record<string, any>; // Step-specific params
    };

    output: {
        response: string;     // AI response
        artifacts: Artifact[]; // Generated files/diagrams
        structuredData?: any;  // Parsed structured output
    };

    // Model & Cost
    modelUsed: string;        // Which AI model
    modelMode: ModelMode;     // thinking | fast | creative | ..
    tokensUsed: number;       // Total tokens
    costCents: number;        // Cost in cents
    latencyMs: number;        // Execution time

    // Quality
    confidence: number;       // 0-1 confidence score
    reasoning: string;        // Explanation of confidence

    // Parallel Execution (if applicable)
    wasParallel: boolean;
    parallelModels?: string[]; // Models used in parallel
    synthesisStrategy?: string; // How results were combined

    // Timestamps
    startedAt: Date;
    completedAt: Date;
```



## 4.2 Session Data Model

### SESSION DATA MODEL

#### SESSION

sessionId: uuid  
tenantId: uuid  
userId: uuid  
title: string  
status: SessionStatus  
domainMode: DomainMode  
createdAt: timestamp  
updatedAt: timestamp

has many

#### CONVERSATIONS

conversationId: uuid  
sessionId: uuid (FK)  
title: string  
createdAt: timestamp

has many

#### MESSAGES

messageId: uuid  
conversationId: uuid (FK)  
role: 'user' | 'assistant' | 'system'  
content: text  
createdAt: timestamp

has many

#### STEPS

stepId: uuid  
sessionId: uuid (FK)  
stepOrder: integer



stepType: StepType  
input: jsonb  
output: jsonb  
modelUsed: string  
tokensUsed: integer  
costCents: decimal  
confidence: decimal  
startedAt: timestamp  
completedAt: timestamp

has many

#### ARTIFACTS

artifactId: uuid  
stepId: uuid (FK)  
type: 'code' | 'diagram' | 'document' | 'data'  
filename: string  
mimeType: string  
s3Key: string  
sizeBytes: integer  
createdAt: timestamp

---

## 5. Collaboration Features

### 5.1 Real-Time Collaboration

#### REAL-TIME COLLABORATION

Think Tank Session  
"Architecture Design #42"

User A  
(Owner)

User B  
(Editor)

User C  
(Viewer)

WebSocket Connection  
(Real-time event streaming)

Event Types

- `step.started`      - A new step is executing
- `step.progress`      - Step progress update
- `step.completed`    - Step finished with result
- `message.added`      - New message in conversation
- `cursor.moved`      - User cursor position
- `user.joined`        - New collaborator joined
- `user.left`         - Collaborator left
- `artifact.created`   - New artifact generated
- `session.status`    - Session state changed

COLLABORATION ROLES:

Role	Permissions
Owner	Full control, manage collaborators, delete session
Editor	Add messages, trigger steps, view all content
Viewer	Read-only access to session and results
Commenter	View + add comments, no step triggering

---

## 6. Domain Modes

### 6.1 Specialized Reasoning Modes

DOMAIN MODES

Think Tank adapts its reasoning approach based on problem domain:

## RESEARCH MODE

Best for: Academic research, literature review, fact-finding

Models: Perplexity Sonar, Claude (deep\_research mode)

Features:

- Source citation
- Cross-reference verification
- Comprehensive literature synthesis

## ENGINEERING MODE

Best for: System design, architecture, technical problems

Models: Claude, GPT-4o, DeepSeek (code mode)

Features:

- Code generation as artifacts
- Architecture diagrams
- Technical trade-off analysis

## ANALYTICAL MODE

Best for: Data analysis, math, statistics, quantitative problems

Models: o1, Claude (thinking mode), DeepSeek R1

Features:

- Step-by-step mathematical reasoning
- Statistical analysis
- Proof verification

## CREATIVE MODE

Best for: Writing, brainstorming, ideation, design

Models: Claude, GPT-4o (creative mode, high temperature)

Features:

- Multiple creative alternatives
- Iterative refinement
- Style adaptation

## LEGAL MODE

Best for: Contract analysis, compliance, legal research

Models: Claude (precise mode), GPT-4o

Features:

- Citation of legal precedents
- Risk assessment
- Compliance checking

#### MEDICAL MODE (HIPAA Compliant)

Best for: Clinical analysis, medical research (non-diagnostic)

Models: Claude (precise mode), approved medical models

Features:

- PHI sanitization
- Medical literature citation
- Disclaimer generation

#### BUSINESS MODE

Best for: Strategy, planning, market analysis, business problems

Models: GPT-4o, Claude, Gemini

Features:

- Framework application (SWOT, Porter's, etc.)
- Financial modeling
- Competitive analysis

#### GENERAL MODE

Best for: Mixed problems, general questions

Models: Automatically selected based on sub-problem analysis

Features:

- Dynamic mode switching per step
- Balanced approach

---

## 7. Quality & Confidence

### 7.1 Confidence Scoring System

#### CONFIDENCE SCORING SYSTEM

Every step and the final solution receives a confidence score (0-1):

#### CONFIDENCE FACTORS

Factor	Contribution
Model Agreement	+0.2 if parallel models agree
Reasoning Depth	+0.15 for thorough chain-of-thought
Source Quality	+0.15 for cited/verified sources
Task Complexity	-0.1 for very complex sub-problems
Model Confidence	+0.1 for high model self-confidence
Consistency	+0.1 for consistency with prior steps
Verification	+0.2 if verified by second model

#### CONFIDENCE LEVELS

0.9 - 1.0	VERY HIGH - Strong consensus
0.7 - 0.9	HIGH - Reliable
0.5 - 0.7	MODERATE - Review recommended
0.3 - 0.5	LOW - Uncertain
0.0 - 0.3	VERY LOW - Needs verification

#### FINAL SOLUTION CONFIDENCE

Formula:

$\text{final\_confidence} = \text{weighted\_avg}(\text{step\_confidences}) \times \text{synthesis\_factor}$

Where:

- step weights based on importance/complexity
- synthesis\_factor accounts for integration quality

8. User Interface

8.1 Think Tank UI Layout

THINK TANK USER INTERFACE		
Think Tank [New] [Share] [Export]		
Logo	Problem: "Design microservices architecture..."	
	Mode: Engineering	Confidence: 0.88 Cost: \$0.18
SESSIONS	MAIN CONVERSATION	DETAILS
Today	STEPS	
Arch #42	You	
Data Q	Design a scalable microservices architecture for an e-commerce platform that handles 10M...	Step 1 0.92
Yesterday		Step 2
ML Model	0.89	
Security		Step 3
	0.91	
Last Week	Think Tank	Step 4
API Des		Running
Budget	I'll approach this problem by:	Step 5
		Step 6
	1. Analyzing requirements...	
	2. Identifying services...	
[+ New]	3. Designing data flow...	ARTIFACTS
	arch.md	
	Step 4 Progress: 65%	diagram
		docker
	Analyzing data patterns...	
		MODELS USED
		Claude 3.5
	GPT-4o	
	Ask a follow-up question...	o1

---

## Think Tank Platform Architecture v3.2.0

*Advanced AI reasoning for complex problems*

---

© 2024 RADIANT. All Rights Reserved.

## AGI & Workflow Orchestration

### Intelligent Multi-Model AI Orchestration

Version 4.18.0 | December 2024

---

#### 1. Overview

RADIANT's AGI Orchestration Layer coordinates multiple AI models using 49 proven patterns to achieve **50-300% quality improvement** over single-model approaches through intelligent model selection, parallel execution, and result synthesis.

#### Why 50-300% Improvement?

RADIANT achieves dramatic quality improvements through four synergistic capabilities:

#### QUALITY IMPROVEMENT ARCHITECTURE

1. FULL MULTI-AI ORCHESTRATION (+50-100%)
  - Multiple models work together on same problem
  - Different models catch different errors
  - Consensus mechanisms eliminate hallucinations
  - Parallel execution = best of all models
  
2. SIMULATED AGI ORCHESTRATION (+75-150%)
  - Intelligent task analysis and decomposition
  - Dynamic model + mode selection per sub-task
  - Self-reflection and metacognition
  - Automatic pattern selection based on problem type
  - Confidence scoring and quality gates

3. SPECIALTY SELF-HOSTED MODELS (56 models) (+50-100%)
  - Domain-specific fine-tuned models
  - Code-specialized models (DeepSeek, CodeLlama)
  - Math-specialized models
  - Legal/Medical/Financial domain models
  - No rate limits, full control
  
4. 49 RESEARCH-BACKED PATTERNS (+25-75%)
  - AI Debate: adversarial improvement
  - Self-Refine: iterative enhancement
  - Chain-of-Verification: fact-checking pipeline
  - Tree of Thoughts: exploration of solution space

COMBINED EFFECT: 50-300% QUALITY IMPROVEMENT

### Key Capabilities

Feature	Description	Improvement
<b>49 Patterns</b>	Proven orchestration workflows from AI research	+25-75%
<b>106+ Models</b>	50 external + 56 self-hosted specialty models	+50-100%
<b>Simulated AGI</b>	Intelligent orchestration with metacognition	+75-150%
<b>9 Modes</b>	Thinking, Research, Fast, Creative, Precise, Code, Vision, Long-context, Standard	+25-50%
<b>Parallel Execution</b>	Multiple models simultaneously with synthesis	+50-100%

### Improvement by Use Case

Use Case	Single Model	RADIANT Orchestrated	Improvement
Complex coding	60% accuracy	95% accuracy	<b>+58%</b>
Legal analysis	70% accuracy	96% accuracy	<b>+37%</b>
Research synthesis	65% completeness	98% completeness	<b>+51%</b>
Creative writing	Good quality	Publication-ready	<b>+100-200%</b>
Multi-step reasoning	55% correct	94% correct	<b>+71%</b>
Fact verification	75% accurate	99% accurate	<b>+32%</b>
Code review	Catches 60% bugs	Catches 95% bugs	<b>+58%</b>



---

## 2. The 49 Orchestration Patterns

### Pattern Categories

#### CATEGORY 1: CONSENSUS & AGGREGATION (Patterns 1-7)

- Self-Consistency (SC)
- Universal Self-Consistency
- Multi-Agent Debate Voting
- Diverse Verifier (DiVeRSe)
- Meta-Reasoning
- Ensemble Refinement
- Sample-and-Marginalize

#### CATEGORY 2: DEBATE & DELIBERATION (Patterns 8-14)

- AI Debate (SOD)
- Multi-Agent Debate
- Consultancy Model
- Society of Mind
- Cross-Examination
- Red-Team/Blue-Team
- Adversarial Collaboration

#### CATEGORY 3: CRITIQUE & REFINEMENT (Patterns 15-21)

- Self-Refine
- Reflexion
- Constitutional AI
- CRITIC
- Recursive Criticism
- Iterative Refinement
- Self-Taught Reasoner

#### CATEGORY 4: VERIFICATION & VALIDATION (Patterns 22-28)

- Chain-of-Verification
- Fact-Checking Pipeline
- Step-by-Step Verification
- Process Reward Model
- Outcome Reward Model
- Dual-Process Verification
- LLM-as-Judge

#### CATEGORY 5: DECOMPOSITION (Patterns 29-35)

- Least-to-Most
- Decomposed Prompting
- Tree of Thoughts
- Skeleton-of-Thought
- Plan-and-Solve

Graph of Thoughts  
Recursive Decomposition

#### CATEGORY 6: SPECIALIZED REASONING (Patterns 36-42)

Chain-of-Thought (CoT)  
ReAct  
Self-Ask  
Maieutic Prompting  
Analogical Reasoning  
Contrastive CoT  
Program-Aided Language Model

#### CATEGORY 7: MULTI-MODEL ROUTING (Patterns 43-46)

Mixture of Experts  
Speculative Decoding  
FrugalGPT  
Model Cascading

#### CATEGORY 8: ENSEMBLE METHODS (Patterns 47-49)

Model Ensemble  
Boosted Prompting  
Blended RAG

---

### 3. AGI Dynamic Model Selection

#### How It Works

##### AGI MODEL SELECTION FLOW

PROMPT: "Write recursive TSP algorithm with dynamic programming"

#### 1. DOMAIN DETECTION

Keywords: "algorithm", "recursive", "programming"  
Detected: CODING (0.85)

#### 2. TASK ANALYSIS

- Complexity: HIGH
- Requires Reasoning: YES
- Requires Precision: YES

```

3. QUERY LIVE MODEL METADATA
modelMetadataService.getAllMetadata()
Returns: 106 models with capabilities, pricing

```

#### 4. SCORE & SELECT WITH MODES

Model	Score	Mode
Claude 3.5 Sonnet	0.94	thinking
OpenAI o1	0.92	thinking
DeepSeek R1	0.88	code

#### Domain Detection Keywords

Domain	Keywords	Best Models
<b>coding</b>	code, function, algorithm, debug	Claude, o1, DeepSeek
<b>math</b>	calculate, equation, proof, theorem	o1, Claude, DeepSeek R1
<b>reasoning research</b>	think, logic, step by step, why comprehensive, investigate, explore	o1, Claude, DeepSeek R1 Perplexity, Gemini Deep
<b>creative</b>	write, story, imagine, design	Claude, GPT-4o

#### 4. Model Execution Modes

Mode	Icon	Auto-Selected When	Parameters
<b>thinking</b>		requiresReasoning + o1/claude/r1	thinkingBudget: 10000
<b>deep_research</b>		requiresResearch + perplexity	searchDepth: comprehensive
<b>fast</b>		flash/turbo/mini models	maxTokens: 2048
<b>creative</b>		requiresCreativity	temperature: 0.9
<b>precise</b>		requiresPrecision	temperature: 0.1

Mode	Icon	Auto-Selected When	Parameters
<b>code</b>		coding domain	temperature: 0.2
<b>vision</b>		vision-capable models	enableVision: true
<b>long_context</b>		large context windows	maxTokens: 16384
<b>standard</b>		default fallback	default params

## 5. Parallel Execution

### Execution Modes

Mode	Behavior	Latency	Best For
<b>all</b>	Wait for all models	Slowest model	Maximum quality
<b>race</b>	First success wins	Fastest model	Low latency
<b>quorum</b>	Wait for X%	Second fastest	Balance

### Synthesis Strategies

Strategy	How It Works
<b>best_of</b>	Select highest confidence response
<b>vote</b>	Choose most common answer (majority)
<b>weighted</b>	Score by confidence $\times$ (1/latency)
<b>merge</b>	AI combines all responses into one

## 6. Visual Workflow Editor

### Editor Features

- **Method Palette** - Drag-and-drop 16 method types
- **Canvas** - Visual workflow with nodes and connections
- **Step Configuration** - 4 tabs: General, Params, Parallel, Advanced
- **Zoom/Pan** - Canvas navigation controls
- **Test & Save** - Execute and persist workflows

### Step Configuration

[General] [Params] [Parallel] [Advanced]

#### PARALLEL TAB

Enable Parallel Execution [ON]  
 AGI Model Selection [ON]

Min Models: [2]      Max Models: [5]  
Domain Hints: [coding, reasoning]

Preferred Modes:  
[ ] thinking   [ ] deep\_research   [ ] fast  
[ ] creative   [ ] precise           [ ] code

Execution Mode: [All (wait for all)]  
Synthesis: [Weighted (confidence + speed)]  
Timeout: [30000] ms

---

## 7. API Usage

### Execute Workflow

```
const result = await orchestrationService.executeWorkflow({
  tenantId: 'tenant-123',
  workflowCode: 'SOD', // AI Debate pattern
  prompt: 'Should we prioritize AI safety over capabilities?',
  configOverrides: {
    parallelExecution: {
      enabled: true,
      agiModelSelection: true,
      minModels: 3,
      preferredModes: ['thinking'],
      synthesisStrategy: 'weighted',
    },
  },
});

// Result includes:
// - response: Final synthesized answer
// - confidence: 0-1 quality score
// - steps: Array of step results
// - modelsUsed: Models that participated
// - totalCost: Cost in cents
// - totalLatency: Time in ms
```

---

## 8. Benefits

Benefit	Single Model	Orchestrated AI
<b>Accuracy</b>	~75%	~92%
<b>Bias</b>	Single perspective	Multi-perspective

Benefit	Single Model	Orchestrated AI
<b>Verification</b>	None	Built-in
<b>Confidence</b>	Unknown	Measured
<b>Reliability</b>	One point of failure	Redundant

---

## RADIANT AGI Orchestration v4.18.0

*Intelligent multi-model AI coordination*

## Simultaneous Prompt Execution

### Overview

Both RADIANT and Think Tank support **simultaneous prompt execution** - the ability to run multiple AI prompts in parallel across different models. This capability enables dramatic quality improvements through consensus mechanisms and significant throughput gains for high-volume applications.

---

## RADIANT Parallel Execution

### Configuration

```
interface ParallelExecutionConfig {
    enabled: boolean;
    mode: 'all' | 'race' | 'quorum';
    models: string[];
    minModels?: number;
    maxModels?: number;
    agiModelSelection?: boolean;
    domainHints?: string[];
    timeoutMs?: number;
    failureStrategy: 'fail_fast' | 'best_effort';
}
```

### Execution Modes

Mode	Description	Use Case
<b>all</b>	Wait for all models to complete	Consensus, synthesis
<b>race</b>	Return first successful response	Speed-critical
<b>quorum</b>	Return when majority agree	Balanced quality/speed

### Implementation

```
export class ParallelExecutionService {
    async executeParallel(
```

```

    prompt: string,
    config: ParallelExecutionConfig
): Promise<ParallelResult> {
    const models = config.agiModelSelection
        ? await this.agiSelectModels(prompt, config)
        : config.models;

    // Launch all models simultaneously
    const promises = models.map(model =>
        this.executeWithTimeout(prompt, model, config.timeoutMs)
    );

    switch (config.mode) {
        case 'all':
            return this.waitForAll(promises);
        case 'race':
            return this.waitForFirst(promises);
        case 'quorum':
            return this.waitForQuorum(promises, config.minModels);
    }
}

private async waitForAll(
    promises: Promise<ModelResponse>[]
): Promise<ParallelResult> {
    const results = await Promise.allSettled(promises);
    const successful = results
        .filter((r): r is PromiseFulfilledResult<ModelResponse> =>
            r.status === 'fulfilled')
        .map(r => r.value);

    // Synthesize consensus from all responses
    const synthesis = await this.synthesizeResponses(successful);

    return {
        responses: successful,
        synthesis,
        consensusScore: this.calculateConsensus(successful),
        totalLatencyMs: Math.max(...successful.map(r => r.latencyMs))
    };
}

private async waitForQuorum(
    promises: Promise<ModelResponse>[],
    minModels: number = Math.ceil(promises.length / 2)
): Promise<ParallelResult> {
    const results: ModelResponse[] = [];

```

```

return new Promise((resolve) => {
  promises.forEach(async (promise) => {
    try {
      const result = await promise;
      results.push(result);

      if (results.length >= minModels) {
        // Check if results agree
        const consensus = this.checkConsensus(results);
        if (consensus.agreement >= 0.7) {
          resolve({
            responses: results,
            synthesis: consensus.synthesized,
            consensusScore: consensus.agreement,
            earlyTermination: true
          });
        }
      }
    } catch (error) {
      // Continue waiting for other models
    }
  });
});
}
}

```

## Use Cases

### 1. Consensus Verification

```

// Run same prompt on 3 models, synthesize agreement
const result = await parallelExecution.executeParallel(
  "What is the capital of France?",
  {
    enabled: true,
    mode: 'all',
    models: ['claude-3-5-sonnet', 'gpt-4o', 'gemini-1.5-pro'],
    agiModelSelection: false
  }
);
// consensusScore: 1.0 - all models agree "Paris"

```

### 2. Code Review with Multiple Perspectives

```

// Different models find different issues
const result = await parallelExecution.executeParallel(
  codeToReview,
  {
    enabled: true,

```



```

    mode: 'all',
    models: ['claude-3-5-sonnet', 'deepseek-coder-v2', 'gpt-4o'],
    domainHints: ['code', 'security', 'performance']
  }
);
// Synthesis combines all found issues

```

### 3. Creative Writing Enhancement

```

// Generate multiple creative variations
const result = await parallelExecution.executeParallel(
  "Write a tagline for an AI company",
  {
    enabled: true,
    mode: 'all',
    models: ['claude-3-5-sonnet', 'gpt-4o'],
    // High temperature for diversity
  }
);
// Get best elements from each response

```

---

## Think Tank Concurrent Sessions

### Session-Level Parallelism

Think Tank supports concurrent execution at multiple levels:

1. **Parallel Steps** - Independent reasoning steps run simultaneously
2. **Multi-Model Steps** - Same step runs on multiple models
3. **Parallel Sessions** - Multiple sessions execute concurrently

### Implementation

```

export class ConcurrentSessionService {
  // Execute independent steps in parallel
  async executeParallelSteps(
    sessionId: string,
    steps: ThinkTankStep[]
  ): Promise<StepResult[]> {
    // Identify which steps can run in parallel
    const { independent, dependent } = this.analyzeDependen(steps);

    // Run independent steps simultaneously
    const independentResults = await Promise.all(
      independent.map(step => this.executeStep(sessionId, step))
    );

    // Run dependent steps sequentially

```

```

    const dependentResults = [];
    for (const step of dependent) {
        dependentResults.push(await this.executeStep(sessionId, step));
    }

    return [...independentResults, ...dependentResults];
}

// Run same step on multiple models for consensus
async executeWithMultipleModels(
    sessionId: string,
    step: ThinkTankStep,
    models: string[]
): Promise<ConsensusResult> {
    // Execute simultaneously on all models
    const responses = await Promise.all(
        models.map(model =>
            this.executeStepWithModel(sessionId, step, model)
        )
    );

    // Synthesize consensus
    return this.synthesizeConsensus(responses);
}

// Parallel problem decomposition
async parallelDecompose(
    sessionId: string,
    problem: string
): Promise<DecompositionResult> {
    // Multiple models decompose the problem differently
    const decompositions = await Promise.all([
        this.decomposeWith(problem, 'claude-3-5-sonnet'),
        this.decomposeWith(problem, 'gpt-4o'),
        this.decomposeWith(problem, 'gemini-1.5-pro')
    ]);

    // Merge decompositions for comprehensive coverage
    return this.mergeDecompositions(decompositions);
}
}

```

## Session Configuration

```

interface ThinkTankSessionConfig {
    sessionId: string;
    parallelExecution: {
        enabled: boolean;
    };
}

```

```

    maxConcurrentSteps: number;           // Default: 5
    maxConcurrentModels: number;          // Default: 3
    consensusThreshold: number;           // 0-1, default: 0.7
    timeoutPerStepMs: number;             // Default: 30000
};
modelSelection: {
    automatic: boolean;                   // AGI selects models
    preferredModels: string[];
    domainHint: string;
};
}

```

## Database Schema Support

```

-- Session configuration for parallel execution
ALTER TABLE thinktank_sessions ADD COLUMN parallel_execution_config JSONB DEFAULT '{
    "enabled": true,
    "maxConcurrentSteps": 5,
    "maxConcurrentModels": 3,
    "consensusThreshold": 0.7
}';

-- Track parallel step executions
CREATE TABLE thinktank_parallel_executions (
    id UUID PRIMARY KEY DEFAULT uuid_generate_v4(),
    session_id UUID NOT NULL REFERENCES thinktank_sessions(id),
    step_id UUID NOT NULL REFERENCES thinktank_steps(id),
    model_id VARCHAR(100) NOT NULL,
    started_at TIMESTAMPTZ NOT NULL,
    completed_at TIMESTAMPTZ,
    response TEXT,
    tokens_used INTEGER,
    latency_ms INTEGER,
    included_in_consensus BOOLEAN DEFAULT true
);

```

## Performance Benefits

### Throughput Improvement

Scenario	Sequential	Parallel	Improvement
3 models, consensus	9s	3.5s	<b>2.6x faster</b>
5-step reasoning	25s	8s	<b>3.1x faster</b>
Code review (3 perspectives)	12s	4.5s	<b>2.7x faster</b>

### Quality Improvement from Consensus

Task	Single Model	3-Model Consensus	Improvement
Fact verification	85%	99%	<b>+16%</b>
Code correctness	78%	95%	<b>+22%</b>
Reasoning accuracy	72%	94%	<b>+31%</b>

## API Usage

### REST API

POST /api/v1/chat/completions  
Content-Type: application/json  
Authorization: Bearer {api\_key}

```
{
  "messages": [{"role": "user", "content": "Explain quantum computing"}],
  "parallel": {
    "enabled": true,
    "mode": "all",
    "models": ["claude-3-5-sonnet", "gpt-4o", "gemini-1.5-pro"]
  }
}
```

### Response

```
{
  "id": "par_abc123",
  "object": "parallel.completion",
  "responses": [
    {"model": "claude-3-5-sonnet", "content": "...", "latency_ms": 2100},
    {"model": "gpt-4o", "content": "...", "latency_ms": 1800},
    {"model": "gemini-1.5-pro", "content": "...", "latency_ms": 2300}
  ],
  "synthesis": {
    "content": "...",
    "consensus_score": 0.92,
    "method": "weighted_merge"
  },
  "usage": {
    "total_tokens": 4521,
    "total_cost_usd": 0.0234
  }
}
```

### SDK Usage

```
import { RadiantClient } from '@radiant/sdk';
```

```

const client = new RadiantClient({ apiKey: 'your-key' });

// Parallel execution
const result = await client.chat.completions.create({
  messages: [{ role: 'user', content: 'Analyze this contract...' }],
  parallel: {
    enabled: true,
    mode: 'all',
    models: ['claude-3-5-sonnet', 'gpt-4o']
  }
});

console.log(result.synthesis.content);
console.log(`Consensus: ${result.synthesis.consensus_score}`);

```

---

## Cost Considerations

Parallel execution uses multiple models, which increases costs but provides:

Trade-off	Single Model	Parallel (3 models)
Cost	\$0.01	\$0.03
Quality	75%	95%
Latency	3s	3.5s
Reliability	99%	99.99%

## Cost-Effective Strategies:

1. Use parallel for critical tasks only
2. Start with cheaper models, escalate if disagreement
3. Use quorum mode to terminate early on agreement
4. Cache consensus results for repeated queries # RADIANT & Think Tank Complete Features List

## Comprehensive Feature Reference

Version 4.18.0 | December 2024

---

## Feature Categories

1. AI Model Management
2. Orchestration & Workflows
3. Think Tank Platform
4. Billing & Cost Management
5. Multi-Tenant Platform
6. Security & Compliance
7. Analytics & Monitoring

- 8. Developer Tools
- 9. Admin Dashboard
- 10. Swift Deployer App

## 1. AI Model Management

### 1.1 Model Router Service

Feature	Description	How It Fits
<b>Unified API</b>	Single API endpoint for 106+ AI models	Developers use one API regardless of provider
<b>Model Fallback</b>	Automatic failover to backup models	Ensures reliability when primary model fails
<b>Rate Limiting</b>	Per-tenant and per-model limits	Prevents abuse and manages costs
<b>Request Routing</b>	Intelligent routing to optimal provider	Minimizes latency, maximizes availability

### 1.2 Model Metadata Service

Feature	Description	How It Fits
<b>Live Model Data</b>	Real-time model availability and capabilities	AGI uses current data for model selection
<b>Capability Scores</b>	0-1 scores for reasoning, coding, creative, etc.	Enables intelligent model matching to tasks
<b>Pricing Data</b>	Input/output token costs per model	Supports cost estimation and budgeting
<b>AI Research</b>	Automated metadata updates via AI	Keeps model info current without manual work
<b>Admin Override</b>	Manual corrections to AI-gathered data	Admins can fix inaccuracies

### 1.3 Supported Models (106+)

Provider	Models	Specialties
<b>OpenAI</b>	GPT-4o, GPT-4o-mini, o1, o1-mini, o3	General, reasoning, multimodal
<b>Anthropic</b>	Claude 3.5 Sonnet, Claude 3 Opus/Haiku	Reasoning, coding, safety
<b>Google</b>	Gemini 2.0 Flash/Pro, Gemini Deep Research	Speed, multimodal, research
<b>Meta</b>	Llama 3.1 (8B/70B/405B)	Open source, customizable
<b>Mistral</b>	Mistral Large, Codestral	European, code

Provider	Models	Specialties
<b>DeepSeek</b>	DeepSeek R1, DeepSeek Chat	Reasoning, cost-effective
<b>Perplexity xAI</b>	Sonar Pro, Sonar	Real-time research
<b>Cohere</b>	Grok 2	Real-time knowledge
<b>+6 more</b>	Command R+, Embed	Enterprise, RAG
	56 self-hosted models	Custom deployments

## 2. Orchestration & Workflows

### 2.1 Orchestration Patterns (49)

Feature	Description	How It Fits
<b>Pattern Library</b>	49 proven multi-AI workflows	Pre-built solutions for complex tasks
<b>Pattern Selection</b>	Automatic best pattern for task	Users don't need to know which pattern to use
<b>Custom Workflows</b>	Create/modify workflow patterns	Tenants can build their own patterns

**Pattern Categories:** - Consensus & Aggregation (7) - Debate & Deliberation (7) - Critique & Refinement (7) - Verification & Validation (7) - Decomposition (7) - Specialized Reasoning (7) - Multi-Model Routing (4) - Ensemble Methods (3)

### 2.2 AGI Dynamic Model Selection

Feature	Description	How It Fits
<b>Domain Detection</b>	Identifies coding, math, legal, etc. from prompt	Matches models to domain expertise
<b>Task Analysis</b>	Detects complexity, reasoning needs	Selects appropriate model count and modes
<b>Live Scoring</b>	Scores all available models for task	Always uses best current models
<b>Mode Assignment</b>	Selects optimal mode per model	Maximizes each model's effectiveness

### 2.3 Model Execution Modes (9)

Mode	Description	How It Fits
<b>Thinking</b>	Extended reasoning (o1, Claude)	Complex problems requiring deep thought
<b>Deep Research</b>	Comprehensive research (Perplexity)	Fact-finding, literature review

Mode	Description	How It Fits
<b>Fast</b>	Speed-optimized (Flash models)	Quick queries, autocomplete
<b>Creative</b>	High temperature output	Writing, brainstorming
<b>Precise</b>	Low temperature, factual	Data extraction, compliance
<b>Code</b>	Code-optimized settings	Programming tasks
<b>Vision</b>	Multimodal with images	Image analysis
<b>Long Context</b>	Extended context window	Large documents
<b>Standard</b>	Default parameters	General use

## 2.4 Parallel Execution

Feature	Description	How It Fits
<b>Multi-Model Calls</b>	Execute 2-10 models simultaneously	Higher quality through diversity
<b>Execution Modes</b>	All, Race, Quorum	Balance quality vs latency
<b>Result Synthesis</b>	Best-of, Vote, Weighted, Merge	Combine multiple responses optimally
<b>Timeout Handling</b>	Per-model timeouts	Prevents slow models from blocking
<b>Failure Strategy</b>	Fail-fast, Continue, Fallback	Graceful degradation

## 2.5 Visual Workflow Editor

Feature	Description	How It Fits
<b>Drag-and-Drop</b>	Visual workflow design	Non-technical users can build workflows
<b>Method Palette</b>	16 reusable method types	Building blocks for any workflow
<b>Step Configuration</b>	4-tab config panel	Fine-grained control per step
<b>Canvas Controls</b>	Zoom, pan, fit	Navigate complex workflows
<b>Test &amp; Save</b>	Execute and persist	Validate before deployment

# 3. Think Tank Platform

## 3.1 Problem Solving Engine

Feature	Description	How It Fits
<b>Problem Decomposition</b>	Breaks complex problems into parts	Makes hard problems tractable
<b>Multi-Step Reasoning</b>	Chain-of-thought with recorded steps	Transparent reasoning process
<b>Solution Synthesis</b>	Combines step outputs into answer	Coherent final solutions



Feature	Description	How It Fits
<b>Confidence Scoring</b>	0-1 quality score per step and overall	Users know reliability

### 3.2 Session Management

Feature	Description	How It Fits
<b>Persistent Sessions</b>	Save and resume any session	Long-running problem solving
<b>Session History</b>	All steps recorded with metadata	Audit trail, learning
<b>Conversation Threads</b>	Multiple conversations per session	Organize follow-ups
<b>Artifact Storage</b>	Code, diagrams, documents as outputs	Tangible deliverables

### 3.3 Domain Modes (8)

Mode	Description	How It Fits
<b>Research</b>	Academic research, fact-finding	Source citation, verification
<b>Engineering</b>	System design, architecture	Code artifacts, diagrams
<b>Analytical</b>	Math, statistics, data analysis	Step-by-step proofs
<b>Creative</b>	Writing, ideation, design	Multiple alternatives
<b>Legal</b>	Contracts, compliance	Risk assessment
<b>Medical</b>	Clinical analysis (HIPAA)	PHI sanitization
<b>Business</b>	Strategy, planning	Framework application
<b>General</b>	Mixed problems	Dynamic mode switching

### 3.4 Collaboration

Feature	Description	How It Fits
<b>Real-Time Sync</b>	WebSocket live updates	Multiple users see changes instantly
<b>Collaboration Roles</b>	Owner, Editor, Viewer, Commenter	Appropriate access control
<b>Cursor Presence</b>	See other users' positions	Awareness of collaborators
<b>Shared Sessions</b>	Invite others to sessions	Team problem solving

## 4. Billing & Cost Management

### 4.1 Credit System

Feature	Description	How It Fits
<b>Credit Accounts</b>	Pre-paid credit balances	Simple usage-based billing
<b>Credit Transactions</b>	Detailed usage history	Transparency on spending
<b>Auto-Refill</b>	Automatic top-up at threshold	Uninterrupted service
<b>Credit Alerts</b>	Low balance notifications	Avoid service interruption

## 4.2 Subscriptions

Feature	Description	How It Fits
<b>Plan Tiers</b>	Free Trial, Individual, Pro, Enterprise	Options for all user types
<b>Feature Gating</b>	Features by plan level	Upsell path
<b>Usage Limits</b>	Tokens/requests per plan	Fair resource allocation
<b>Stripe Integration</b>	Payment processing	Industry-standard payments

## 4.3 Cost Management

Feature	Description	How It Fits
<b>Budget Alerts</b>	Spending limit notifications	Prevent cost overruns
<b>Cost Estimation</b>	Pre-request cost estimates	Informed decisions
<b>Usage Analytics</b>	Spend by model, user, time	Optimize usage patterns
<b>Invoice Generation</b>	Automated monthly invoices	Accounting integration

# 5. Multi-Tenant Platform

## 5.1 Tenant Management

Feature	Description	How It Fits
<b>Tenant Isolation</b>	Complete data separation	Security, privacy
<b>Tenant Settings</b>	Per-tenant configuration	Customization
<b>Tenant Onboarding</b>	Self-service signup	Scalable growth
<b>Tenant Suspension</b>	Disable/enable tenants	Account management

## 5.2 User Management

Feature	Description	How It Fits
<b>User Accounts</b>	Individual user identities	Personalization, audit
<b>Role-Based Access</b>	Admin, User, Viewer roles	Appropriate permissions
<b>User Preferences</b>	Model preferences, settings	Personal customization
<b>User Activity</b>	Usage tracking per user	Analytics, billing

### 5.3 API Key Management

Feature	Description	How It Fits
<b>API Key Generation</b>	Create scoped keys	Programmatic access
<b>Key Rotation</b>	Scheduled key rotation	Security best practice
<b>Key Scopes</b>	Limit key permissions	Least privilege
<b>Key Analytics</b>	Usage per key	Monitor applications

---

## 6. Security & Compliance

### 6.1 Data Security

Feature	Description	How It Fits
<b>Row-Level Security</b>	PostgreSQL RLS policies	Automatic tenant isolation
<b>Encryption at Rest</b>	AES-256 encryption	Data protection
<b>Encryption in Transit</b>	TLS 1.3	Secure communication
<b>KMS Key Management</b>	AWS KMS for secrets	Secure key storage

### 6.2 Authentication

Feature	Description	How It Fits
<b>Cognito Integration</b>	AWS Cognito user pools	Enterprise-grade auth
<b>JWT Tokens</b>	Secure session tokens	Stateless auth
<b>MFA Support</b>	Multi-factor authentication	Enhanced security
<b>SSO/SAML</b>	Enterprise SSO integration	Corporate identity

### 6.3 Compliance

Feature	Description	How It Fits
<b>SOC2 Controls</b>	Security controls	Enterprise compliance
<b>HIPAA Mode</b>	Healthcare compliance	Medical use cases
<b>PHI Sanitization</b>	Automatic PII detection	Protect patient data
<b>Audit Logging</b>	Comprehensive audit trail	Compliance reporting
<b>Data Residency</b>	Region-specific deployment	Regulatory requirements

---

## 7. Analytics & Monitoring

### 7.1 Usage Analytics

Feature	Description	How It Fits
<b>Request Metrics</b>	Requests by model, user, time	Usage patterns
<b>Token Tracking</b>	Input/output token counts	Cost attribution
<b>Latency Metrics</b>	Response time tracking	Performance monitoring
<b>Error Rates</b>	Failure tracking	Reliability monitoring

## 7.2 Model Performance

Feature	Description	How It Fits
<b>Quality Scores</b>	Model quality over time	Identify degradation
<b>Comparison Reports</b>	Model vs model analysis	Model selection
<b>A/B Testing</b>	Test model variations	Optimize choices
<b>Learning Data</b>	ML training data collection	Continuous improvement

## 7.3 Business Intelligence

Feature	Description	How It Fits
<b>Dashboard</b>	Executive metrics view	Quick status
<b>Custom Reports</b>	Build custom analytics	Specific insights
<b>Export</b>	CSV/PDF export	External analysis
<b>Alerts</b>	Threshold notifications	Proactive monitoring

# 8. Developer Tools

## 8.1 SDK

Feature	Description	How It Fits
<b>TypeScript SDK</b>	Type-safe client library	Developer productivity
<b>API Documentation</b>	OpenAPI/Swagger docs	Self-service integration
<b>Code Examples</b>	Sample implementations	Quick start
<b>Playground</b>	Interactive API testing	Experimentation

## 8.2 Webhooks

Feature	Description	How It Fits
<b>Event Webhooks</b>	Push notifications for events	Real-time integrations
<b>Webhook Management</b>	Create, update, delete hooks	Self-service config
<b>Retry Logic</b>	Automatic retry on failure	Reliability
<b>Webhook Logs</b>	Delivery history	Debugging

## 8.3 Integrations

Feature	Description	How It Fits
<b>Slack Integration</b>	Notifications to Slack	Team communication
<b>Zapier Connect</b>	5000+ app integrations	Automation
<b>Custom Webhooks</b>	HTTP POST to any endpoint	Flexible integration

## 9. Admin Dashboard

### 9.1 Dashboard Pages

Page	Description	How It Fits
<b>Overview</b>	System health, key metrics	At-a-glance status
<b>Tenants</b>	Tenant management	Customer administration
<b>Users</b>	User administration	Access control
<b>Models</b>	Model configuration	AI management
<b>Orchestration</b>	Workflow patterns	Pattern management
<b>Analytics</b>	Usage reports	Business intelligence
<b>Billing</b>	Revenue, invoices	Financial management
<b>Security</b>	Audit logs, compliance	Security oversight
<b>Settings</b>	Platform configuration	System settings

### 9.2 UI Features

Feature	Description	How It Fits
<b>Responsive Design</b>	Mobile-friendly	Access anywhere
<b>Dark Mode</b>	Light/dark themes	User preference
<b>Search</b>	Global search	Find anything quickly
<b>Filters</b>	Advanced filtering	Narrow results
<b>Bulk Actions</b>	Multi-select operations	Efficiency

## 10. Swift Deployer App

### 10.1 Deployment Features

Feature	Description	How It Fits
<b>CDK Deployment</b>	One-click AWS deployment	Simple infrastructure setup
<b>Progress Tracking</b>	Real-time deployment status	Visibility into process
<b>Stack Management</b>	Deploy individual stacks	Granular control
<b>Rollback</b>	Revert failed deployments	Safety net

## 10.2 QA & Testing

Feature	Description	How It Fits
<b>Test Suites</b>	Run unit/integration tests	Quality assurance
<b>Test Results</b>	Pass/fail reporting	Quick feedback
<b>Coverage Reports</b>	Code coverage metrics	Quality metrics

## 10.3 AI Assistant

Feature	Description	How It Fits
<b>Deployment Guidance</b>	AI helps with deployment	Reduces errors
<b>Error Diagnosis</b>	AI analyzes failures	Faster resolution
<b>Best Practices</b>	AI suggests improvements	Optimization

## 10.4 Local Storage

Feature	Description	How It Fits
<b>SQLCIPHER DB</b>	Encrypted local storage	Secure credentials
<b>AWS Profiles</b>	Multiple AWS accounts	Environment management
<b>Deployment History</b>	Past deployment records	Audit trail

---

## RADIANT Feature Reference v4.18.0

*106+ models • 49 patterns • 9 modes • Enterprise-grade*

## RADIANT Services Reference

### Complete Lambda Services Inventory (62 Services)

#### Core Infrastructure Services

**1. BrainRouter (brain-router.ts)** **Purpose:** Central routing service that directs incoming requests to appropriate handlers based on task type.

**Key Methods:** - `routeTask(task: Task): Promise<TaskResult>` - Routes task to handler - `analyzeTaskType(input: string): TaskType` - Determines task classification - `selectHandler(taskType: TaskType): Handler` - Selects appropriate handler

**Task Types:** | Type | Description | Handler | |——|———|———| | **generation** | Text generation | ModelRouterService | | **analysis** | Data analysis | AnalyticsService | | **transformation** | Content transformation | TransformService | | **orchestration** | Multi-step workflow | OrchestrationService | | **conversation** | Chat interaction | ConversationService |

**2. ThermalStateService (thermal-state.ts)** **Purpose:** Monitors system thermal state and adjusts workload distribution.

**Key Methods:** - `getSystemState(): ThermalState` - Current system state - `adjustWorkload(state: ThermalState): void` - Modify processing - `recordMetric(name: string, value: number): void` - Track metrics

**States:** - `nominal` - Normal operation - `elevated` - Increased load - `throttled` - Reduced capacity - `critical` - Emergency mode

---

**3. MetricsCollector (metrics-collector.ts)** **Purpose:** Collects and aggregates system metrics for monitoring.

**Key Methods:** - `recordLatency(service: string, ms: number): void` - `recordTokenUsage(model: string, input: number, output: number): void` - `recordCost(tenantId: string, cents: number): void` - `getMetrics(timeRange: TimeRange): MetricsSummary`

**Metrics Tracked:** - API latency (p50, p95, p99) - Token usage by model - Cost by tenant - Error rates - Provider health

---

**4. ErrorLogger (error-logger.ts)** **Purpose:** Structured error logging with context preservation.

**Key Methods:** - `logError(error: Error, context: ErrorContext): void` - `logWarning(message: string, data: object): void` - `getRecentErrors(count: number): ErrorLog[]`

**Error Categories:** - `PROVIDER_ERROR` - AI provider failures - `VALIDATION_ERROR` - Input validation - `AUTH_ERROR` - Authentication failures - `RATE_LIMIT` - Rate limiting triggered - `INTERNAL_ERROR` - System errors

---

**5. CredentialsManager (credentials-manager.ts)** **Purpose:** Secure management of API keys and credentials.

**Key Methods:** - `getCredential(provider: string): Promise<string>` - `rotateCredential(provider: string): Promise<void>` - `validateCredential(provider: string): Promise<boolean>`

**Supported Providers:** - OpenAI, Anthropic, Google, Mistral - Groq, Perplexity, xAI, Together - Cohere, DeepSeek, Replicate

---

## AI Model Services

**6. ModelRouterService (model-router.service.ts)** **Purpose:** Routes AI requests to optimal provider with fallback.

### Architecture:

Request → Validate → Select Provider → Execute → Fallback (if needed) → Response

**Model Registry (24 Models):**

Model ID	Provider	Capabilities	Cost (\$/1K tokens)
anthropic/claude-3-5-sonnet	bedrock	reasoning, coding, vision	\$0.003/\$0.015
anthropic/claude-3-haiku	bedrock	fast, efficient	\$0.00025/\$0.00125
meta/llama-3.1-70b	bedrock	reasoning, open-source	\$0.00099/\$0.00099
amazon/titan-text-express	bedrock	fast, aws-native	\$0.0002/\$0.0006
openai/gpt-4o	litellm	reasoning, vision	\$0.005/\$0.015
openai/gpt-4o-mini	litellm	fast, efficient	\$0.00015/\$0.0006
openai/o1	litellm	reasoning, math	\$0.015/\$0.060
openai/o1-mini	litellm	reasoning, coding	\$0.003/\$0.012
google/gemini-1.5-pro	litellm	reasoning, long-context	\$0.00125/\$0.005
google/gemini-1.5-flash	litellm	fast, vision	\$0.000075/\$0.0003
mistral/mistral-large	litellm	reasoning, multilingual	\$0.003/\$0.009
mistral/codestral	litellm	coding	\$0.001/\$0.003
cohere/command-r-plus	litellm	reasoning, rag	\$0.003/\$0.015
deepseek/deepseek-coder-v2	litellm	coding	\$0.00014/\$0.00028
groq/llama-3.1-70b-versatile	groq	fast, reasoning	\$0.00059/\$0.00079
groq/llama-3.1-8b-instant	groq	instant, fast	\$0.00005/\$0.00008
groq/mixtral-8x7b	groq	fast, moe	\$0.00024/\$0.00024
perplexity/sonar-large	perplexity	search, citations	\$0.001/\$0.001
perplexity/sonar-small	perplexity	search, fast	\$0.0002/\$0.0002
xai/grok-beta	xai	reasoning, realtime	\$0.005/\$0.015
together/llama-3.1-405b	together	reasoning, large	\$0.005/\$0.015

**Fallback Chains:**

```

bedrock → litellm → groq
litellm → bedrock → groq
groq → litellm → bedrock
perplexity → litellm

```



xai → litellm → groq  
together → litellm → groq

**Provider Health Tracking:** - isHealthy: boolean - latencyMs: number - errorCount: number  
- consecutiveFailures: number ( $\geq 3$  marks unhealthy)

---

**7. ModelMetadataService (model-metadata.service.ts)** **Purpose:** Manages live model capabilities, pricing, and availability.

**Key Methods:** - getMetadata(modelId: string): Promise<ModelMetadata> - getAllMetadata(): Promise<ModelMetadata[]> - updateMetadata(modelId: string, data: Partial<ModelMetadata>): Promise<void> - refreshFromInternet(): Promise<void> - AI-powered metadata updates

**Metadata Structure:**

```
interface ModelMetadata {
  modelId: string;
  provider: string;
  displayName: string;
  description: string;
  capabilities: {
    reasoning: number;      // 0-1 score
    coding: number;
    creative: number;
    factual: number;
    math: number;
    vision: boolean;
    longContext: boolean;
  };
  contextWindow: number;
  maxOutputTokens: number;
  pricing: {
    inputPer1kTokens: number;
    outputPer1kTokens: number;
    currency: string;
  };
  availability: {
    isAvailable: boolean;
    regions: string[];
    lastChecked: Date;
  };
  performance: {
    avgLatencyMs: number;
    throughputTokensPerSec: number;
  };
}
```

---

**8. ModelSelectionService (model-selection-service.ts)** **Purpose:** Intelligent model selection based on task characteristics.

**Selection Algorithm:** 1. **Domain Detection** - Identify problem domain from keywords 2. **Task Analysis** - Determine complexity, requirements 3. **Model Scoring** - Score each model for task fit 4. **Mode Assignment** - Select optimal execution mode 5. **Cost/Quality Balance** - Apply user preferences

**Domain Keywords:** | Domain | Keywords | |——|———| | coding | code, function, algorithm, debug, implement, API | | math | calculate, equation, formula, solve, proof | | legal | contract, law, compliance, regulation, liability | | medical | diagnosis, treatment, symptom, clinical, patient | | research | study, analyze, evidence, literature, methodology | | creative | write, story, design, brainstorm, creative |

---

## Orchestration Services

**9. OrchestrationPatternsService (orchestration-patterns.service.ts)** **Purpose:** Manages 49 orchestration patterns with parameterized methods.

### Pattern Categories (8):

Category	Count	Examples
Consensus & Aggregation	7	Self-Consistency, Meta-Reasoning, Mixture-of-Agents
Debate & Deliberation	7	AI Debate, Society of Mind, Socratic Dialogue
Critique & Refinement	7	Self-Refine, Reflexion, Constitutional AI
Verification & Validation	7	Chain-of-Verification, LLM-as-Judge, Fact-Check
Decomposition	7	Least-to-Most, Tree of Thoughts, Skeleton-of-Thought
Specialized Reasoning	7	Chain-of-Thought, ReAct, Self-Ask
Multi-Model Routing	4	Mixture of Experts, FrugalGPT, Cascading
Ensemble Methods	3	Model Ensemble, Blended RAG, Speculative Decoding

### All 49 Patterns:

1. **Self-Consistency** - Multiple samples, majority vote
2. **Universal Self-Consistency** - Free-form answer selection
3. **Meta-Reasoning** - Compare reasoning paths
4. **DiVeRSe** - Diverse verifier ensemble
5. **Mixture-of-Agents** - Multi-agent aggregation
6. **LLM-Blender** - Pairwise ranking fusion

7. **Multi-Agent Consensus** - Agent negotiation
8. **AI Debate** - Adversarial debate with judge
9. **Multi-Agent Debate** - Multi-party debate
10. **Society of Mind** - Agent specialization
11. **ChatEval** - Multi-agent evaluation
12. **ReConcile** - Confidence-weighted discussion
13. **Socratic Dialogue** - Question-based exploration
14. **Diplomatic Consensus** - Negotiated agreement
15. **Self-Refine** - Iterative refinement
16. **Reflexion** - Verbal reinforcement learning
17. **CRITIC** - External tool verification
18. **Iterative Refinement** - Multi-pass improvement
19. **Constitutional AI** - Principle-based critique
20. **Progressive Refinement** - Staged quality improvement
21. **Expert Refinement** - Domain expert review
22. **Chain-of-Verification** - Claim verification chain
23. **LLM-as-Judge** - Model evaluation
24. **Self-Verification** - Self-checking
25. **G-Eval** - Structured evaluation
26. **Cross-Validation** - Multi-model validation
27. **Fact-Check Chain** - Fact verification pipeline
28. **Consensus Validation** - Agreement-based validation
29. **Least-to-Most** - Simple to complex decomposition
30. **Decomposed Prompting** - Sub-task breakdown
31. **Tree of Thoughts** - Branching exploration
32. **Graph of Thoughts** - Graph-based reasoning
33. **Skeleton-of-Thought** - Parallel point expansion
34. **Plan-and-Solve** - Planning then execution
35. **Recursive Decomposition** - Hierarchical breakdown
36. **Chain-of-Thought** - Step-by-step reasoning
37. **Self-Ask** - Sub-question generation
38. **ReAct** - Reasoning + Acting
39. **Program-of-Thoughts** - Code-based reasoning
40. **Analogical Reasoning** - Example-based reasoning
41. **Maieutic Prompting** - Tree explanation
42. **Contrastive CoT** - Valid/invalid contrast
43. **Mixture of Experts** - Specialized routing
44. **FrugalGPT** - Cost-optimized cascading
45. **Router Chain** - Capability-based routing
46. **Speculative Routing** - Predictive routing
47. **Model Ensemble** - Multi-model combination
48. **Blended RAG** - RAG ensemble
49. **Speculative Decoding** - Draft-verify acceleration

---

**10. WorkflowEngine (workflow-engine.ts)** **Purpose:** Executes DAG-based workflows with task dependencies.

**Key Methods:** - `createWorkflow(definition: WorkflowDefinition): Promise<string>` - `addTask(workflowId: string, task: Task): Promise<void>` - `startExecution(workflowId: string, params: object): Promise<string>` - `updateExecutionStatus(executionId: string, status: Status): Promise<void>`

### Workflow Definition:

```
interface WorkflowDefinition {
  workflowId: string;
  name: string;
  description: string;
  category: 'generation' | 'analysis' | 'transformation' | 'pipeline' | 'custom';
  dagDefinition: {
    nodes: TaskNode[];
    edges: Edge[];
  };
  inputSchema: JSONSchema;
  outputSchema: JSONSchema;
  defaultParameters: Record<string, any>;
  timeoutSeconds: number;
  maxRetries: number;
}

interface TaskNode {
  taskId: string;
  taskType: 'model_inference' | 'transformation' | 'condition' | 'parallel' | 'aggregation';
  config: object;
  dependsOn: string[];
  conditionExpression?: string;
}
```

---

**11. ResponseSynthesisService (response-synthesis.service.ts)** **Purpose:** Synthesizes responses from multiple AI models.

### Synthesis Strategies:

Strategy	Description	Best For
best_of	Select highest confidence response	Quality-critical
vote	Majority voting on answer	Factual questions
weighted	Confidence $\times$ (1/latency) weighted	Balanced
merge	AI combines all responses	Complex analysis

### Merge Algorithm:

1. Collect all responses with metadata
2. Extract key points from each
3. Identify agreements and conflicts

4. Generate unified response
5. Apply conflict resolution
6. Calculate final confidence

---

## Billing Services

**12. BillingService (billing.ts)** **Purpose:** Manages credits, subscriptions, and billing.

**Key Methods:** - `getSubscription(tenantId: string): Promise<Subscription>` - `getCreditBalance(tenantId: string): Promise<CreditBalance>` - `addCredits(tenantId: string, amount: number, type: string): Promise<number>` - `useCredits(tenantId: string, amount: number): Promise<{success, newBalance}>` - `purchaseCredits(tenantId: string, amount: number, price: number): Promise<string>`

**Subscription Tiers:** | Tier | Monthly Price | Annual Price | Credits/User | |——|———|———|  
 ——|———| | Free Trial | \$0 | - | 100 | | Individual | \$19 | \$190 | 1,000 | | Pro | \$49 | \$490 | 5,000 | | Team | \$199 | \$1,990 | 25,000 | | Enterprise | Custom | Custom | Custom |

**Volume Discounts:** | Credit Amount | Discount | Bonus Credits | |———|———|———|  
 —| | 10-19 | 5% | 0 | | 20-49 | 10% | 0 | | 50-99 | 15% | 5% | | 100+ | 25% | 10% |

**Transaction Types:** - `purchase` - Credit purchase - `bonus` - Promotional credits - `refund` - Refunded credits - `usage` - Credits consumed - `transfer_in` / `transfer_out` - Credit transfers - `subscription_allocation` - Monthly allocation - `expiration` - Expired credits - `adjustment` - Manual adjustment

---

**13. StorageBillingService (storage-billing.ts)** **Purpose:** Tracks storage costs per tenant.

**Billable Storage:** - Uploaded files - Generated artifacts - Session history - Conversation logs

**Pricing:** - \$0.023 per GB/month (Standard) - \$0.0125 per GB/month (Infrequent) - \$0.004 per GB/month (Archive)

---

## Cognitive Services

**14. CognitiveBrainService (cognitive-brain.service.ts)** **Purpose:** High-level cognitive processing and reasoning.

**Cognitive Capabilities:** - Working memory management - Attention allocation - Abstract reasoning - Analogy formation - Concept learning

---

**15. ReasoningEngine (reasoning-engine.ts)** **Purpose:** Chain-of-thought and multi-step reasoning.

**Reasoning Modes:** | Mode | Description | |——|—————| | **deductive** | From general to specific  
| | **inductive** | From specific to general | | **abductive** | Best explanation inference | | **analogical**  
| Similarity-based reasoning |

---

**16. CausalReasoningService (causal-reasoning.service.ts)** **Purpose:** Causal inference and counterfactual reasoning.

**Methods:** - `identifyCauses(effect: string): Promise<Cause[]>` - `predictEffects(cause: string): Promise<Effect[]>` - `counterfactual(scenario: string, change: string): Promise<string>`

---

**17. GoalPlanningService (goal-planning.service.ts)** **Purpose:** Goal decomposition and planning.

**Planning Algorithm:**

1. Parse high-level goal
  2. Identify subgoals
  3. Determine dependencies
  4. Sequence actions
  5. Allocate resources
  6. Execute and monitor
- 

**18. MetacognitionService (metacognition.service.ts)** **Purpose:** Self-reflection and learning from mistakes.

**Metacognitive Functions:** - Confidence calibration - Error detection - Strategy selection - Performance monitoring

---

## Memory Services

**19. MemoryService (memory-service.ts)** **Purpose:** Persistent memory across sessions.

**Memory Types:** - **Short-term:** Current session context - **Long-term:** Cross-session knowledge  
- **Episodic:** Event-based memories - **Semantic:** Factual knowledge

---

**20. EpisodicMemoryService (episodic-memory.service.ts)** **Purpose:** Event-based memory storage and retrieval.

**Key Methods:** - `recordEpisode(event: Episode): Promise<void>` - `retrieveRelevant(query: string, limit: number): Promise<Episode[]>` - `consolidate(): Promise<void>` - Memory optimization

---

**21. MemoryConsolidationService (memory-consolidation.service.ts)** **Purpose:** Optimizes memory storage by consolidating similar memories.

---

**22. TimeMachineService (time-machine.ts)** **Purpose:** Access historical state at any point in time.

**Key Methods:** - `getStateAt(timestamp: Date): Promise<SystemState>` - `getDiff(from: Date, to: Date): Promise<StateDiff>` - `restore(timestamp: Date): Promise<void>`

---

## AGI Services

**23. AGIOrchestratorService (agi-orchestrator.service.ts)** **Purpose:** Coordinates AGI capabilities across services.

---

**24. AdvancedAGIService (advanced-agi.service.ts)** **Purpose:** Advanced AGI features including self-improvement.

---

**25. AGICompleteService (agi-complete.service.ts)** **Purpose:** Complete AGI pipeline from input to output.

---

**26. AGIExtensionsService (agi-extensions.service.ts)** **Purpose:** Extensible AGI capabilities.

---

## Collaboration Services

**27. CollaborationService (collaboration.ts)** **Purpose:** Real-time collaboration features.

**WebSocket Events:** | Event | Direction | Description | |——-|———|———-| | `join_session` | Client→Server | Join collaborative session | | `leave_session` | Client→Server | Leave session | | `cursor_move` | Bidirectional | Cursor position update | | `content_update` | Bidirectional | Content change | | `user_joined` | Server→Client | New user notification | | `user_left` | Server→Client | User left notification |

---

**28. ConcurrentSessionManager (concurrent-session.ts)** **Purpose:** Manages concurrent user sessions.

**Key Methods:** - `createSession(config: SessionConfig): Promise<string>` - `joinSession(sessionId: string, userId: string): Promise<void>` - `getSessionState(sessionId: string):`

```
Promise<SessionState>    -    broadcastUpdate(sessionId: string, update: Update):
Promise<void>
```

---

**29. TeamService (team-service.ts)** **Purpose:** Team and organization management.

**Key Methods:** - createTeam(tenantId: string, name: string): Promise<string> -  
 addMember(teamId: string, userId: string, role: string): Promise<void>-getTeamMembers(teamId:  
 string): Promise<Member[]>

---

## Additional Services (30-62)

#	Service	File	Purpose
30	NeuralEngine	neural-engine.ts	Neural network operations
31	AutoResolveService	auto-resolve.ts	Automatic conflict resolution
32	CanvasService	canvas-service.ts	Visual canvas artifacts
33	PersonaService	persona-service.ts	AI persona management
34	SchedulerService	scheduler-service.ts	Task scheduling
35	LicenseService	license-service.ts	License management
36	UnifiedModelRegistry	unified-model-registry.ts	Unified model registry
37	GrandfatheringService	grandfathering-service.ts	Legacy migration
38	VoiceVideoService	voice-video.ts	Voice/video processing
39	ResultMergingService	result-merging.ts	Merge results
40	WorldModelService	world-model.service.ts	World state modeling
41	MultiAgentService	multi-agent.service.ts	Multi-agent coordination
42	TheoryOfMindService	theory-of-mind.service.ts	Mental state modeling
43	MultimodalBindingService	multimodal-binding.service.ts	Cross-modal binding
44	SkillExecutionService	skill-execution.service.ts	Skill execution
45	AutonomousAgentService	autonomous-agent.service.ts	Autonomous operations
46	ConsciousnessService	consciousness.service.ts	Consciousness modeling
47	ConfigEngineService	config-engine.service.ts	Configuration engine
48	SelfImprovementService	self-improvement.service.ts	Self-improvement
49	MoralCompassService	moral-compass.service.ts	Ethical reasoning
50	MLTrainingService	ml-training.service.ts	ML model training
51	LearningService	learning.service.ts	Learning data collection
52	FeedbackService	feedback.service.ts	User feedback
53	FeedbackLearningService	feedback-learning.ts	Learn from feedback
54	WorkflowProposalService	workflow-proposals.ts	Workflow improvements
55	AppIsolationService	app-isolation.ts	App-level isolation
56	LocalizationService	localization.ts	i18n support
57	MigrationApprovalService	migration-approval.ts	Migration approval
58	SuperiorOrchestrationService	superior-orchestration.ts	Superior services
59	RadiantUnifiedService	radiant-unified.service.ts	Unified API
60	NeuralOrchestrationService	neural-orchestration.ts	Neural orchestration
61	AuditService	audit.ts	Audit logging



#	Service	File	Purpose
62	APIKeysService	api-keys.ts	API key management

## RADIANT Database Schema Reference

### Complete Migration Inventory (40 Migrations)

#### Migration Index

#	Migration	Tables Created	Purpose
001	initial_schema	tenants, users, administrators, invitations, approval_requests	Core platform tables
002	tenant_isolation	RLS policies	Row-level security
003	ai_models	ai_models, model_capabilities	Model registry
004	usage_billing	usage_records, invoices	Usage tracking
005	admin_approval	approval_workflows	Admin approvals
006	self_hosted_models	self_hosted_models, model_deployments	Self-hosted AI
007	external_providers	external_providers, provider_configs	Provider management
010	visual_ai_pipeline	visual_pipelines, pipeline_stages	Visual AI processing
011	brain_router	routing_rules, task_classifications	Request routing
012	metrics_analytics	metrics, analytics_snapshots	Analytics
013	neural_engine	neural_configs, neural_executions	Neural processing
014	error_logging	error_logs, error_patterns	Error tracking
015	credentials_registry	credentials, credential_rotations	Credential management
016	think_tank	thinktank_sessions, thinktank_steps, thinktank_tools	Think Tank
017	concurrent_chat	chat_sessions, chat_messages	Chat management
018	realtime_collaboration	collaborations, collaboration_members	Real-time collab
019	persistent_memory	memories, memory_associations	Memory system
020	focus_personas	personas, persona_configs	AI personas
021	team_plans	teams, team_members, team_plans	Team management
022	provider_registry	providers, provider_health	Provider registry
023	time_machine	snapshots, snapshot_diffs	Time machine
024	orchestration_engine	workflow_definitions, workflow_tasks, workflow_executions, task_executions	Orchestration
025	license_management	licenses, license_activations	Licensing
026	unified_model_registry	unified_models, model_versions	Model registry
027	feedback_learning	feedback, learning_samples	Feedback system
028	neural_orchestration	neural_workflows, neural_steps	Neural orchestration

#	Migration	Tables Created	Purpose
029	workflow_proposals	proposals, proposal_evidence	Workflow improvements
030	app_isolation	app_contexts, app_permissions	App isolation
031	internationalization	locales, translations	i18n
032	dynamic_configuration	configs, config_history	Dynamic config
033	billing_credits	credit_balances, credit_transactions, credit_purchases	Credits
034	storage_billing	storage_usage, storage_costs	Storage billing
035	versioned_subscriptions	subscription_tiers, subscriptions	Subscriptions
036	dual_admin_approval	dual_approvals, approval_chains	Dual approval
037	canvas_artifacts	canvases, canvas_elements	Canvas
038	scheduled_prompts	scheduled_prompts, prompt_executions	Scheduling
039	auto_resolve	auto_resolutions, resolution_rules	Auto-resolve
040	model_selection_pricing	model_pricing, selection_history	Pricing
041	admin_billing_enhancements	billing_reports, revenue_tracking	Billing enhancements

## Core Tables (Migration 001)

### tenants

Primary multi-tenant organization table.

```
CREATE TABLE tenants (
  id UUID PRIMARY KEY DEFAULT uuid_generate_v4(),
  name VARCHAR(100) NOT NULL UNIQUE,
  display_name VARCHAR(200) NOT NULL,
  domain VARCHAR(255),
  settings JSONB NOT NULL DEFAULT '{}',
  status VARCHAR(20) NOT NULL DEFAULT 'active'
    CHECK (status IN ('active', 'suspended', 'pending')),
  created_at TIMESTAMPTZ NOT NULL DEFAULT NOW(),
  updated_at TIMESTAMPTZ NOT NULL DEFAULT NOW()
);

-- Indexes
CREATE INDEX idx_tenants_name ON tenants(name);
CREATE INDEX idx_tenants_domain ON tenants(domain) WHERE domain IS NOT NULL;
CREATE INDEX idx_tenants_status ON tenants(status);
```

Settings JSON Structure:

```
{
  "branding": {
    "logo_url": "string",
```

```

    "primary_color": "#hex",
    "company_name": "string"
  },
  "limits": {
    "max_users": 100,
    "max_api_keys": 10,
    "monthly_token_limit": 1000000
  },
  "features": {
    "think_tank_enabled": true,
    "orchestration_enabled": true,
    "collaboration_enabled": true
  },
  "compliance": {
    "hipaa_mode": false,
    "data_retention_days": 90
  }
}

```

---

## users

End users within tenants.

```

CREATE TABLE users (
  id UUID PRIMARY KEY DEFAULT uuid_generate_v4(),
  tenant_id UUID NOT NULL REFERENCES tenants(id) ON DELETE CASCADE,
  cognito_user_id VARCHAR(128) NOT NULL,
  email VARCHAR(255) NOT NULL,
  display_name VARCHAR(200),
  role VARCHAR(50) NOT NULL DEFAULT 'user'
    CHECK (role IN ('user', 'power_user', 'admin')),
  status VARCHAR(20) NOT NULL DEFAULT 'active'
    CHECK (status IN ('active', 'suspended', 'pending')),
  settings JSONB NOT NULL DEFAULT '{}',
  created_at TIMESTAMPTZ NOT NULL DEFAULT NOW(),
  updated_at TIMESTAMPTZ NOT NULL DEFAULT NOW(),
  UNIQUE (tenant_id, email),
  UNIQUE (cognito_user_id)
);

```

*-- Indexes*

```

CREATE INDEX idx_users_tenant_id ON users(tenant_id);
CREATE INDEX idx_users_email ON users(email);
CREATE INDEX idx_users_cognito_user_id ON users(cognito_user_id);
CREATE INDEX idx_users_status ON users(status);

```

**User Roles:** | Role | Permissions | |——|—————| | user | Basic API access, own resources | |

power\_user | + Create API keys, advanced features | | admin | + Manage users, view analytics |

---

## administrators

Platform administrators (separate from tenant users).

```
CREATE TABLE administrators (  
  id UUID PRIMARY KEY DEFAULT uuid_generate_v4(),  
  cognito_user_id VARCHAR(128) NOT NULL UNIQUE,  
  email VARCHAR(255) NOT NULL UNIQUE,  
  display_name VARCHAR(200) NOT NULL,  
  role VARCHAR(50) NOT NULL DEFAULT 'admin'  
    CHECK (role IN ('super_admin', 'admin', 'operator', 'auditor')),  
  permissions TEXT[] NOT NULL DEFAULT '{}',  
  mfa_enabled BOOLEAN NOT NULL DEFAULT false,  
  last_login_at TIMESTAMPTZ,  
  created_at TIMESTAMPTZ NOT NULL DEFAULT NOW(),  
  updated_at TIMESTAMPTZ NOT NULL DEFAULT NOW(),  
  invited_by UUID REFERENCES administrators(id)  
);
```

**Admin Roles:** | Role | Description | Permissions | |——|———|———| | super\_admin |  
Full platform access | All operations | | admin | Standard admin | Manage tenants, users, models |  
| operator | Operations | View logs, manage deployments | | auditor | Read-only audit | View all,  
modify none |

---

## approval\_requests

Two-person approval system.

```
CREATE TABLE approval_requests (  
  id UUID PRIMARY KEY DEFAULT uuid_generate_v4(),  
  requester_id UUID NOT NULL REFERENCES administrators(id),  
  action_type VARCHAR(100) NOT NULL,  
  resource_type VARCHAR(100) NOT NULL,  
  resource_id VARCHAR(255),  
  payload JSONB NOT NULL DEFAULT '{}',  
  status VARCHAR(20) NOT NULL DEFAULT 'pending'  
    CHECK (status IN ('pending', 'approved', 'rejected', 'expired')),  
  required_approvals INTEGER NOT NULL DEFAULT 1,  
  approvals JSONB NOT NULL DEFAULT '[]',  
  expires_at TIMESTAMPTZ NOT NULL,  
  created_at TIMESTAMPTZ NOT NULL DEFAULT NOW(),  
  updated_at TIMESTAMPTZ NOT NULL DEFAULT NOW()  
);
```

**Action Types Requiring Approval:** - delete\_tenant - Delete a tenant - modify\_billing

- Change billing settings - grant\_super\_admin - Elevate to super admin - bulk\_data\_export - Export all data - disable\_security\_feature - Disable security

---

## Think Tank Tables (Migration 016)

### thinktank\_sessions

Problem-solving session tracking.

```
CREATE TABLE thinktank_sessions (  
  id UUID PRIMARY KEY DEFAULT uuid_generate_v4(),  
  tenant_id UUID NOT NULL REFERENCES tenants(id) ON DELETE CASCADE,  
  user_id UUID NOT NULL REFERENCES users(id) ON DELETE CASCADE,  
  problem_summary TEXT,  
  domain VARCHAR(50),  
  complexity VARCHAR(20) CHECK (complexity IN ('low', 'medium', 'high', 'extreme')),  
  total_steps INTEGER DEFAULT 0,  
  avg_confidence DECIMAL(3, 2),  
  solution_found BOOLEAN DEFAULT false,  
  total_tokens INTEGER DEFAULT 0,  
  total_cost DECIMAL(10, 6) DEFAULT 0,  
  created_at TIMESTAMPTZ NOT NULL DEFAULT NOW(),  
  completed_at TIMESTAMPTZ  
);  
  
-- Indexes  
CREATE INDEX idx_thinktank_sessions_tenant ON thinktank_sessions(tenant_id, created_at DESC);  
CREATE INDEX idx_thinktank_sessions_user ON thinktank_sessions(user_id);  
  
-- RLS  
ALTER TABLE thinktank_sessions ENABLE ROW LEVEL SECURITY;  
CREATE POLICY thinktank_sessions_isolation ON thinktank_sessions  
  FOR ALL USING (tenant_id = current_setting('app.current_tenant_id', true)::uuid);
```

**Domains:** - research - Academic research - engineering - Technical problems - analytical - Data analysis - creative - Creative tasks - legal - Legal analysis - medical - Medical queries (HIPAA) - business - Business strategy - general - General problems

---

### thinktank\_steps

Individual reasoning steps within sessions.

```
CREATE TABLE thinktank_steps (  
  id UUID PRIMARY KEY DEFAULT uuid_generate_v4(),  
  session_id UUID NOT NULL REFERENCES thinktank_sessions(id) ON DELETE CASCADE,  
  step_number INTEGER NOT NULL,  
  step_type VARCHAR(50) NOT NULL
```

```

        CHECK (step_type IN ('decompose', 'reason', 'execute', 'verify', 'synthesize')),
description TEXT,
reasoning TEXT,
result TEXT,
confidence DECIMAL(3, 2) CHECK (confidence >= 0 AND confidence <= 1),
model_used VARCHAR(100),
tokens_used INTEGER,
duration_ms INTEGER,
created_at TIMESTAMPTZ NOT NULL DEFAULT NOW()
);

-- Indexes
CREATE INDEX idx_thinktank_steps_session ON thinktank_steps(session_id, step_number);

-- RLS
ALTER TABLE thinktank_steps ENABLE ROW LEVEL SECURITY;
CREATE POLICY thinktank_steps_isolation ON thinktank_steps
    FOR ALL USING (
        session_id IN (SELECT id FROM thinktank_sessions
                        WHERE tenant_id = current_setting('app.current_tenant_id', true)::uuid)
    );

```

**Step Types:**

Type	Description	Typical Model	Step	Model
decompose	Break problem into parts	Claude 3.5	reason	Chain-of-thought reasoning
execute	Execute solution step	Task-specific	verify	Verify result accuracy
synthesize	Combine into final answer	Claude 3.5		

## thinktank\_tools

Available tools for Think Tank.

```

CREATE TABLE thinktank_tools (
    id UUID PRIMARY KEY DEFAULT uuid_generate_v4(),
    tool_name VARCHAR(100) NOT NULL UNIQUE,
    tool_type VARCHAR(50) NOT NULL,
    description TEXT,
    parameters_schema JSONB NOT NULL DEFAULT '{}',
    implementation TEXT,
    is_active BOOLEAN DEFAULT true,
    created_at TIMESTAMPTZ NOT NULL DEFAULT NOW()
);

-- Default tools
INSERT INTO thinktank_tools (tool_name, tool_type, description, parameters_schema) VALUES
    ('web_search', 'search', 'Search the web for information', '{"query": "string"}'),
    ('calculator', 'compute', 'Perform mathematical calculations', '{"expression": "string"}'),
    ('code_executor', 'compute', 'Execute code snippets', '{"language": "string", "code": "string"}');

```

```
('file_reader', 'io', 'Read file contents', '{"path": "string"}'),
('api_caller', 'network', 'Make API requests', '{"url": "string", "method": "string", "body": "string"}'),
```

---

## Orchestration Tables (Migration 024)

### workflow\_definitions

Workflow pattern definitions.

```
CREATE TABLE workflow_definitions (
  id UUID PRIMARY KEY DEFAULT uuid_generate_v4(),
  workflow_id VARCHAR(100) NOT NULL UNIQUE,
  name VARCHAR(200) NOT NULL,
  description TEXT,
  category VARCHAR(50) NOT NULL
    CHECK (category IN ('generation', 'analysis', 'transformation', 'pipeline', 'custom')),
  version VARCHAR(20) NOT NULL DEFAULT '1.0.0',

  dag_definition JSONB NOT NULL DEFAULT '{}',
  input_schema JSONB NOT NULL DEFAULT '{}',
  output_schema JSONB NOT NULL DEFAULT '{}',
  default_parameters JSONB NOT NULL DEFAULT '{}',

  timeout_seconds INTEGER DEFAULT 3600,
  max_retries INTEGER DEFAULT 3,
  min_tier INTEGER DEFAULT 1,

  is_active BOOLEAN DEFAULT true,
  requires_audit_trail BOOLEAN DEFAULT false,
  hipaa_compliant BOOLEAN DEFAULT false,

  created_at TIMESTAMPTZ NOT NULL DEFAULT NOW(),
  updated_at TIMESTAMPTZ NOT NULL DEFAULT NOW(),
  created_by UUID
);
```

### DAG Definition Structure:

```
{
  "nodes": [
    {
      "taskId": "decompose",
      "taskType": "model_inference",
      "modelId": "anthropic/claude-3-5-sonnet",
      "config": {
        "systemPrompt": "Break down the problem...",
        "temperature": 0.3
      }
    },
  ],
}
```

```

    "dependsOn": []
  },
  {
    "taskId": "solve_part_1",
    "taskType": "model_inference",
    "dependsOn": ["decompose"]
  }
],
"edges": [
  {"from": "decompose", "to": "solve_part_1"}
]
}

```

---

## workflow\_tasks

Individual tasks within workflows.

```

CREATE TABLE workflow_tasks (
  id UUID PRIMARY KEY DEFAULT uuid_generate_v4(),
  workflow_id UUID NOT NULL REFERENCES workflow_definitions(id) ON DELETE CASCADE,
  task_id VARCHAR(100) NOT NULL,
  name VARCHAR(200) NOT NULL,
  description TEXT,

  task_type VARCHAR(50) NOT NULL
    CHECK (task_type IN ('model_inference', 'transformation', 'condition',
                        'parallel', 'aggregation', 'external_api', 'human_review')),
  model_id VARCHAR(100),
  service_id VARCHAR(100),

  config JSONB NOT NULL DEFAULT '{}',
  input_mapping JSONB DEFAULT '{}',
  output_mapping JSONB DEFAULT '{}',

  sequence_order INTEGER DEFAULT 0,
  depends_on TEXT[] DEFAULT '{}',
  condition_expression TEXT,
  timeout_seconds INTEGER DEFAULT 300,

  created_at TIMESTAMPTZ NOT NULL DEFAULT NOW(),
  updated_at TIMESTAMPTZ NOT NULL DEFAULT NOW(),
  UNIQUE(workflow_id, task_id)
);

```

**Task Types:**

Type	Description	Example
model_inference	AI model call	Generate text
transformation	Data transformation	Format output
condition	Conditional branching	If confidence > 0.8
parallel	Parallel execution	Call 3 models



| aggregation | Combine results | Merge responses | | external\_api | External API call | Web search | | human\_review | Human-in-the-loop | Approval step |

---

## workflow\_executions

Workflow execution tracking.

```
CREATE TABLE workflow_executions (  
  id UUID PRIMARY KEY DEFAULT uuid_generate_v4(),  
  workflow_id UUID NOT NULL REFERENCES workflow_definitions(id),  
  tenant_id UUID NOT NULL REFERENCES tenants(id) ON DELETE CASCADE,  
  user_id UUID NOT NULL,  
  
  status VARCHAR(20) NOT NULL DEFAULT 'pending'  
    CHECK (status IN ('pending', 'running', 'paused', 'completed', 'failed', 'cancelled'))  
  
  input_parameters JSONB NOT NULL DEFAULT '{}',  
  resolved_parameters JSONB DEFAULT '{}',  
  output_data JSONB,  
  
  error_message TEXT,  
  error_details JSONB,  
  
  started_at TIMESTAMPTZ,  
  completed_at TIMESTAMPTZ,  
  duration_ms INTEGER,  
  
  estimated_cost_usd DECIMAL(10, 4),  
  actual_cost_usd DECIMAL(10, 4),  
  
  checkpoint_data JSONB,  
  priority INTEGER DEFAULT 5 CHECK (priority >= 1 AND priority <= 10),  
  
  created_at TIMESTAMPTZ NOT NULL DEFAULT NOW(),  
  updated_at TIMESTAMPTZ NOT NULL DEFAULT NOW()  
);  
  
-- RLS  
ALTER TABLE workflow_executions ENABLE ROW LEVEL SECURITY;  
CREATE POLICY workflow_executions_isolation ON workflow_executions  
  FOR ALL USING (tenant_id = current_setting('app.current_tenant_id', true)::uuid);
```

---

## task\_executions

Individual task execution tracking.

```

CREATE TABLE task_executions (
  id UUID PRIMARY KEY DEFAULT uuid_generate_v4(),
  workflow_execution_id UUID NOT NULL REFERENCES workflow_executions(id) ON DELETE CASCADE,
  task_id VARCHAR(100) NOT NULL,

  status VARCHAR(20) NOT NULL DEFAULT 'pending'
    CHECK (status IN ('pending', 'running', 'completed', 'failed', 'skipped', 'retrying')),
  attempt_number INTEGER DEFAULT 1,

  input_data JSONB,
  output_data JSONB,

  error_message TEXT,
  error_code VARCHAR(50),

  started_at TIMESTAMPTZ,
  completed_at TIMESTAMPTZ,
  duration_ms INTEGER,

  resource_usage JSONB DEFAULT '{}',
  cost_usd DECIMAL(10, 4),

  created_at TIMESTAMPTZ NOT NULL DEFAULT NOW()
);

```

---

## Billing Tables (Migrations 033-035)

### credit\_balances

Tenant credit balance tracking.

```

CREATE TABLE credit_balances (
  tenant_id UUID PRIMARY KEY REFERENCES tenants(id) ON DELETE CASCADE,
  balance DECIMAL(12, 4) NOT NULL DEFAULT 0,
  lifetime_purchased DECIMAL(12, 4) NOT NULL DEFAULT 0,
  lifetime_used DECIMAL(12, 4) NOT NULL DEFAULT 0,
  lifetime_bonus DECIMAL(12, 4) NOT NULL DEFAULT 0,
  low_balance_alert_threshold DECIMAL(12, 4),
  last_low_balance_alert TIMESTAMPTZ,
  auto_purchase_enabled BOOLEAN DEFAULT false,
  auto_purchase_threshold DECIMAL(12, 4),
  auto_purchase_amount DECIMAL(12, 4),
  created_at TIMESTAMPTZ NOT NULL DEFAULT NOW(),
  updated_at TIMESTAMPTZ NOT NULL DEFAULT NOW()
);

```

---

## credit\_transactions

Credit transaction history.

```
CREATE TABLE credit_transactions (  
  id UUID PRIMARY KEY DEFAULT uuid_generate_v4(),  
  tenant_id UUID NOT NULL REFERENCES tenants(id) ON DELETE CASCADE,  
  transaction_type VARCHAR(30) NOT NULL  
    CHECK (transaction_type IN ('purchase', 'bonus', 'refund', 'usage',  
                                'transfer_in', 'transfer_out',  
                                'subscription_allocation', 'expiration', 'adjustment')),  
  amount DECIMAL(12, 4) NOT NULL,  
  balance_after DECIMAL(12, 4) NOT NULL,  
  description TEXT,  
  reference_id VARCHAR(255),  
  created_at TIMESTAMPTZ NOT NULL DEFAULT NOW()  
);  
  
-- Indexes  
CREATE INDEX idx_credit_transactions_tenant ON credit_transactions(tenant_id, created_at DESC)
```

---

## subscription\_tiers

Available subscription plans.

```
CREATE TABLE subscription_tiers (  
  id VARCHAR(50) PRIMARY KEY,  
  display_name VARCHAR(100) NOT NULL,  
  description TEXT,  
  price_monthly DECIMAL(10, 2),  
  price_annual DECIMAL(10, 2),  
  included_credits_per_user DECIMAL(10, 2) NOT NULL DEFAULT 0,  
  features JSONB NOT NULL DEFAULT '{}',  
  limits JSONB NOT NULL DEFAULT '{}',  
  is_public BOOLEAN DEFAULT true,  
  sort_order INTEGER DEFAULT 0,  
  created_at TIMESTAMPTZ NOT NULL DEFAULT NOW()  
);  
  
-- Default tiers  
INSERT INTO subscription_tiers (id, display_name, price_monthly, price_annual, included_credits_per_user, features, limits, is_public, sort_order, created_at)  
(  
  ('free', 'Free', 0, NULL, 100, '{"think_tank": false, "orchestration": false, "models": ["gpt-4"]}', '{"tokens": 1000000000, "models": ["gpt-4"]}', true, 0, NOW()),  
  ('pro', 'Pro', 49, 490, 5000, '{"think_tank": true, "orchestration": true, "models": "all"}', '{"tokens": 1000000000, "models": "all"}', true, 1, NOW()),  
  ('team', 'Team', 199, 1990, 25000, '{"think_tank": true, "orchestration": true, "collaboration": true}', '{"tokens": 1000000000, "models": "all"}', true, 2, NOW()),  
  ('enterprise', 'Enterprise', NULL, NULL, 0, '{"think_tank": true, "orchestration": true, "collaboration": true}', '{"tokens": 1000000000, "models": "all"}', true, 3, NOW())  
)
```

---

## subscriptions

Active tenant subscriptions.

```
CREATE TABLE subscriptions (  
    id UUID PRIMARY KEY DEFAULT uuid_generate_v4(),  
    tenant_id UUID NOT NULL REFERENCES tenants(id) ON DELETE CASCADE,  
    tier_id VARCHAR(50) NOT NULL REFERENCES subscription_tiers(id),  
    status VARCHAR(20) NOT NULL DEFAULT 'active'  
        CHECK (status IN ('active', 'cancelled', 'past_due', 'trialing', 'paused')),  
    billing_cycle VARCHAR(10) NOT NULL CHECK (billing_cycle IN ('monthly', 'annual')),  
    seats_purchased INTEGER NOT NULL DEFAULT 1,  
    seats_used INTEGER NOT NULL DEFAULT 0,  
    current_period_start TIMESTAMPTZ NOT NULL,  
    current_period_end TIMESTAMPTZ NOT NULL,  
    cancel_at_period_end BOOLEAN DEFAULT false,  
    cancelled_at TIMESTAMPTZ,  
    stripe_customer_id VARCHAR(255),  
    stripe_subscription_id VARCHAR(255),  
    created_at TIMESTAMPTZ NOT NULL DEFAULT NOW(),  
    updated_at TIMESTAMPTZ NOT NULL DEFAULT NOW()  
);
```

---

## Row-Level Security (Migration 002)

### RLS Pattern

All tenant-scoped tables use this pattern:

```
-- Enable RLS on table  
ALTER TABLE {table_name} ENABLE ROW LEVEL SECURITY;  
  
-- Create isolation policy  
CREATE POLICY {table_name}_isolation ON {table_name}  
    FOR ALL USING (tenant_id = current_setting('app.current_tenant_id', true)::uuid);
```

### Setting Tenant Context

Every request sets the tenant context before queries:

```
SET app.current_tenant_id = '{tenant_uuid}';
```

### Tables with RLS Enabled:

- users
- thinktank\_sessions
- thinktank\_steps
- workflow\_executions
- task\_executions
- credit\_transactions

- subscriptions
  - chat\_sessions
  - chat\_messages
  - collaborations
  - memories
  - feedback
  - api\_keys
  - (all tenant-scoped tables)
- 

## Common Patterns

### Updated At Trigger

```
CREATE OR REPLACE FUNCTION update_updated_at_column()
RETURNS TRIGGER AS $$
BEGIN
    NEW.updated_at = NOW();
    RETURN NEW;
END;
$$ LANGUAGE plpgsql;

-- Apply to tables
CREATE TRIGGER update_{table}_updated_at
    BEFORE UPDATE ON {table}
    FOR EACH ROW EXECUTE FUNCTION update_updated_at_column();
```

### Soft Delete Pattern

```
-- Add columns
deleted_at TIMESTAMPTZ,
deleted_by UUID,

-- Query with filter
WHERE deleted_at IS NULL
```

### Audit Columns

```
created_at TIMESTAMPTZ NOT NULL DEFAULT NOW(),
updated_at TIMESTAMPTZ NOT NULL DEFAULT NOW(),
created_by UUID,
updated_by UUID
```

# Swift Deployer App Reference

## App Architecture

### Overview

The Swift Deployer is a native macOS application for deploying and managing RADIANT infrastructure to AWS.

**Requirements:** - macOS 13.0 (Ventura) or later - Swift 5.9+ - Xcode 15+

### File Structure (36 Files)

```
apps/swift-deployer/  
  Package.swift  
  Sources/RadiantDeployer/  
    RadiantDeployerApp.swift      # App entry point  
    AppState.swift                # Global state management  
  
  Config/  
    RadiantConfig.swift           # App configuration  
  
  Models/ (6 files)  
    Configuration.swift           # Deployment configuration  
    Credentials.swift             # AWS credentials model  
    Deployment.swift              # Deployment state model  
    DomainConfiguration.swift     # Domain settings  
    InstallationParameters.swift  # Install params  
    ManagedApp.swift              # Managed app model  
  
  Services/ (21 files)  
    AIAssistantService.swift      # AI deployment assistant  
    AIRegistryService.swift       # AI model registry  
    APIService.swift              # API communication  
    AWSService.swift              # AWS SDK wrapper  
    AuditLogger.swift             # Audit logging  
    CDKService.swift              # CDK deployment  
    CredentialService.swift       # Credential management  
    DNSService.swift              # DNS configuration  
    DatabaseService.swift         # Local SQLite/SQLCipher  
    DeploymentLockService.swift   # Deployment locking  
    DeploymentService.swift       # Main deployment logic  
    GitHubPackageRegistry.swift   # Package downloads  
    HealthCheckService.swift      # Health monitoring  
    LocalStorageManager.swift     # Encrypted local storage  
    MultiRegionService.swift      # Multi-region deployment  
    OnePasswordService.swift      # 1Password integration  
    PackageService.swift          # Package management  
    SeedDataService.swift         # Database seeding  
    SnapshotService.swift         # State snapshots
```

TimeoutService.swift	# Timeout handling
VoiceInputService.swift	# Voice commands
Views/ (8+ files)	
ABTestingView.swift	# A/B testing config
ContentView.swift	# Main content
DeploymentView.swift	# Deployment UI
SettingsView.swift	# App settings
... (other views)	
Components/ (4 files)	
MacOSComponents.swift	# Design tokens & components
AppCommands.swift	# Menu bar commands
DataTableComponents.swift	# Table components
DetailViewComponents.swift	# Detail view patterns

---

## Models

### Configuration.swift

Deployment configuration model.

```

struct DeploymentConfiguration: Codable, Sendable {
    var appId: String
    var environment: Environment
    var tier: Int
    var region: AWSRegion
    var domain: String?
    var enabledStacks: Set<StackName>
    var customParameters: [String: String]
}

enum Environment: String, Codable, CaseIterable, Sendable {
    case development = "dev"
    case staging = "staging"
    case production = "prod"
}

enum StackName: String, Codable, CaseIterable, Sendable {
    case networking
    case foundation
    case data
    case storage
    case auth
    case ai
    case api
    case admin

```

```

    case batch
    case collaboration
    case monitoring
    case security
    case webhooks
    case scheduledTasks
    case multiRegion
}

```

## Credentials.swift

AWS credentials model.

```

struct AWSCredentials: Codable, Sendable {
    let accessKeyId: String
    let secretAccessKey: String
    let sessionToken: String?
    let region: String
    let profile: String?
    let expiresAt: Date?

    var isExpired: Bool {
        guard let expiresAt else { return false }
        return Date() > expiresAt
    }
}

struct CredentialProfile: Codable, Identifiable, Sendable {
    let id: UUID
    var name: String
    var credentials: AWSCredentials
    var isDefault: Bool
    var lastUsed: Date?
}

```

## Deployment.swift

Deployment state tracking.

```

struct Deployment: Codable, Identifiable, Sendable {
    let id: UUID
    var configuration: DeploymentConfiguration
    var status: DeploymentStatus
    var stackStatuses: [StackName: StackStatus]
    var startedAt: Date
    var completedAt: Date?
    var errorMessage: String?
    var outputs: [String: String]
}

```



```

enum DeploymentStatus: String, Codable, Sendable {
    case pending
    case preparing
    case deploying
    case verifying
    case completed
    case failed
    case rollingBack
    case cancelled
}

enum StackStatus: String, Codable, Sendable {
    case pending
    case creating
    case updating
    case complete
    case failed
    case rollbackInProgress
    case rollbackComplete
    case deleted
}

```

---

## Services

### CDKService.swift

AWS CDK deployment service.

```

actor CDKService {
    private let shell: ShellService
    private let logger: AuditLogger

    // Deploy a single stack
    func deployStack(
        _ stack: StackName,
        config: DeploymentConfiguration,
        credentials: AWSCredentials
    ) async throws -> StackOutput {
        let command = buildCDKCommand(
            action: "deploy",
            stack: stack,
            config: config
        )

        return try await shell.execute(
            command,

```

```

        environment: credentials.asEnvironment()
    )
}

// Deploy all stacks in dependency order
func deployAllStacks(
    config: DeploymentConfiguration,
    credentials: AWSCredentials,
    progressHandler: @Sendable (StackName, StackStatus) -> Void
) async throws -> DeploymentResult {
    let orderedStacks = topologicalSort(config.enabledStacks)
    var outputs: [StackName: StackOutput] = [:]

    for stack in orderedStacks {
        progressHandler(stack, .creating)
        do {
            outputs[stack] = try await deployStack(stack, config: config, credentials: credentials, progressHandler: progressHandler)
            progressHandler(stack, .complete)
        } catch {
            progressHandler(stack, .failed)
            throw DeploymentError.stackFailed(stack, error)
        }
    }

    return DeploymentResult(outputs: outputs)
}

// Stack dependency order
private func topologicalSort(_ stacks: Set<StackName>) -> [StackName] {
    // networking → foundation → data/storage/auth → ai → api/admin → rest
    let order: [StackName] = [
        .networking, .foundation, .data, .storage, .auth,
        .ai, .api, .admin, .batch, .collaboration,
        .monitoring, .security, .webhooks, .scheduledTasks, .multiRegion
    ]
    return order.filter { stacks.contains($0) }
}
}

```

## DeploymentService.swift

Main deployment orchestration.

```

@MainActor
class DeploymentService: ObservableObject {
    @Published var currentDeployment: Deployment?
    @Published var deploymentHistory: [Deployment] = []
    @Published var isDeploying = false
}

```

```

private let cdkService: CDKService
private let healthService: HealthCheckService
private let auditLogger: AuditLogger
private let localStorage: LocalStorageManager

func startDeployment(config: DeploymentConfiguration) async throws {
    guard !isDeploying else {
        throw DeploymentError.alreadyInProgress
    }

    isDeploying = true
    let deployment = Deployment(
        id: UUID(),
        configuration: config,
        status: .preparing,
        stackStatuses: [:],
        startedAt: Date()
    )
    currentDeployment = deployment

    do {
        // 1. Validate configuration
        try await validateConfiguration(config)

        // 2. Acquire deployment lock
        try await acquireLock(config.appId)

        // 3. Deploy stacks
        currentDeployment?.status = .deploying
        let result = try await cdkService.deployAllStacks(
            config: config,
            credentials: try await getCredentials(),
            progressHandler: { [weak self] stack, status in
                Task { @MainActor in
                    self?.currentDeployment?.stackStatuses[stack] = status
                }
            }
        )

        // 4. Verify deployment
        currentDeployment?.status = .verifying
        try await healthService.verifyDeployment(result)

        // 5. Complete
        currentDeployment?.status = .completed
        currentDeployment?.completedAt = Date()
        currentDeployment?.outputs = result.flatOutputs
    }
}

```

```

    } catch {
        currentDeployment?.status = .failed
        currentDeployment?.errorMessage = error.localizedDescription
        throw error
    } finally {
        isDeploying = false
        if let deployment = currentDeployment {
            deploymentHistory.append(deployment)
            try? await localStorage.saveDeployment(deployment)
        }
    }
}

func rollback(deploymentId: UUID) async throws {
    // Rollback implementation
}
}

```

## AIAssistantService.swift

AI-powered deployment assistant.

```

actor AIAssistantService {
    private let apiService: APIService

    struct AssistantResponse: Sendable {
        let message: String
        let suggestions: [Suggestion]
        let actions: [SuggestedAction]
    }

    enum SuggestedAction: Sendable {
        case deployStack(StackName)
        case checkHealth
        case viewLogs(String)
        case runDiagnostics
        case contactSupport
    }

    // Get deployment guidance
    func getDeploymentGuidance(
        for config: DeploymentConfiguration,
        currentStatus: DeploymentStatus?
    ) async throws -> AssistantResponse {
        let prompt = buildGuidancePrompt(config: config, status: currentStatus)
        let response = try await apiService.chat(
            messages: [.init(role: .user, content: prompt)],

```

```

        model: "anthropic/claude-3-haiku"
    )
    return parseAssistantResponse(response)
}

// Diagnose deployment error
func diagnoseError(
    error: Error,
    stackName: StackName,
    logs: String
) async throws -> AssistantResponse {
    let prompt = """
    Analyze this AWS CDK deployment error and provide:
    1. Root cause analysis
    2. Specific fix steps
    3. Prevention recommendations

    Stack: \$(stackName.rawValue)
    Error: \$(error.localizedDescription)
    Logs:
    \$(logs.prefix(2000))
    """

    let response = try await apiService.chat(
        messages: [.init(role: .user, content: prompt)],
        model: "anthropic/claude-3-5-sonnet"
    )
    return parseAssistantResponse(response)
}
}

```

## LocalStorageManager.swift

Encrypted local storage using SQLCipher.

```

actor LocalStorageManager {
    private let db: Connection
    private let encryptionKey: String

    init() throws {
        let path = LocalStorageManager.databasePath
        db = try Connection(path)
        encryptionKey = try KeychainService.getOrCreateDatabaseKey()
        try db.key(encryptionKey)
        try createTablesIfNeeded()
    }

    // Tables

```

```

private func createTablesIfNeeded() throws {
    try db.execute("""
        CREATE TABLE IF NOT EXISTS deployments (
            id TEXT PRIMARY KEY,
            data BLOB NOT NULL,
            created_at INTEGER NOT NULL
        );

        CREATE TABLE IF NOT EXISTS credentials (
            id TEXT PRIMARY KEY,
            profile_name TEXT NOT NULL,
            encrypted_data BLOB NOT NULL,
            is_default INTEGER DEFAULT 0,
            last_used INTEGER
        );

        CREATE TABLE IF NOT EXISTS settings (
            key TEXT PRIMARY KEY,
            value TEXT NOT NULL
        );
    """)
}

// Save deployment
func saveDeployment(_ deployment: Deployment) throws {
    let data = try JSONEncoder().encode(deployment)
    try db.run("""
        INSERT OR REPLACE INTO deployments (id, data, created_at)
        VALUES (?, ?, ?)
    """, deployment.id.uuidString, data, Date().timeIntervalSince1970)
}

// Get deployment history
func getDeploymentHistory() throws -> [Deployment] {
    let rows = try db.prepare("""
        SELECT data FROM deployments ORDER BY created_at DESC LIMIT 100
    """)
    return try rows.compactMap { row in
        guard let data = row[0] as? Data else { return nil }
        return try JSONDecoder().decode(Deployment.self, from: data)
    }
}

// Credential management
func saveCredentials(_ credentials: CredentialProfile) throws {
    let encrypted = try encrypt(credentials)
    try db.run("""
        INSERT OR REPLACE INTO credentials (id, profile_name, encrypted_data, is_default,

```

```

VALUES (?, ?, ?, ?, ?)
""", credentials.id.uuidString, credentials.name, encrypted,
credentials.isDefault ? 1 : 0, credentials.lastUsed?.timeIntervalSince1970)
}
}

```

## HealthCheckService.swift

Deployment health verification.

```

actor HealthCheckService {
    struct HealthCheckResult: Sendable {
        let service: String
        let status: HealthStatus
        let latencyMs: Int?
        let message: String?
    }

    enum HealthStatus: Sendable {
        case healthy
        case degraded
        case unhealthy
        case unknown
    }

    func verifyDeployment(_ result: DeploymentResult) async throws {
        var checks: [HealthCheckResult] = []

        // Check API Gateway
        if let apiUrl = result.outputs["ApiUrl"] {
            checks.append(await checkEndpoint(apiUrl + "/health", service: "API Gateway"))
        }

        // Check LiteLLM
        if let litellmUrl = result.outputs["LiteLLMUrl"] {
            checks.append(await checkEndpoint(litellmUrl + "/health", service: "LiteLLM"))
        }

        // Check Database
        checks.append(await checkDatabase(result.outputs["DatabaseEndpoint"]))

        // Evaluate results
        let unhealthy = checks.filter { $0.status == .unhealthy }
        if !unhealthy.isEmpty {
            throw HealthCheckError.servicesUnhealthy(unhealthy)
        }
    }
}

```

```

private func checkEndpoint(_ url: String, service: String) async -> HealthCheckResult {
    let start = Date()
    do {
        let (_, response) = try await URLSession.shared.data(from: URL(string: url)!)
        let httpResponse = response as! HTTPURLResponse
        let latency = Int(Date().timeIntervalSince(start) * 1000)

        return HealthCheckResult(
            service: service,
            status: httpResponse.statusCode == 200 ? .healthy : .degraded,
            latencyMs: latency,
            message: nil
        )
    } catch {
        return HealthCheckResult(
            service: service,
            status: .unhealthy,
            latencyMs: nil,
            message: error.localizedDescription
        )
    }
}
}

```

---

## Components

### MacOSComponents.swift

Design tokens and reusable components.

```

// Design Tokens
enum RadiantSpacing {
    static let xxs: CGFloat = 2
    static let xs: CGFloat = 4
    static let sm: CGFloat = 8
    static let md: CGFloat = 12
    static let lg: CGFloat = 16
    static let xl: CGFloat = 24
}

enum RadiantRadius {
    static let sm: CGFloat = 4
    static let md: CGFloat = 8
    static let lg: CGFloat = 12
    static let xl: CGFloat = 16
}

```



```

// Status Badge Component
struct StatusBadge: View {
    let status: String
    let color: Color

    var body: some View {
        Text(status)
            .font(.caption)
            .fontWeight(.medium)
            .padding(.horizontal, RadiantSpacing.sm)
            .padding(.vertical, RadiantSpacing.xxs)
            .background(color.opacity(0.15))
            .foregroundColor(color)
            .cornerRadius(RadiantRadius.sm)
    }
}

// Progress Indicator
struct DeploymentProgressView: View {
    let stacks: [StackName]
    let statuses: [StackName: StackStatus]

    var body: some View {
        VStack(alignment: .leading, spacing: RadiantSpacing.sm) {
            ForEach(stacks, id: \.self) { stack in
                HStack {
                    statusIcon(for: statuses[stack] ?? .pending)
                    Text(stack.rawValue)
                        .font(.system(.body, design: .monospaced))
                    Spacer()
                    StatusBadge(
                        status: (statuses[stack] ?? .pending).rawValue,
                        color: statusColor(statuses[stack] ?? .pending)
                    )
                }
            }
        }
    }
}

```

## AppCommands.swift

Menu bar commands with keyboard shortcuts.

```

struct AppCommands: Commands {
    @ObservedObject var appState: AppState

    var body: some Commands {

```

```

CommandGroup(replacing: .newItem) {
    Button("New Deployment") {
        appState.showNewDeploymentSheet = true
    }
    .keyboardShortcut("n", modifiers: .command)

    Button("Import Configuration...") {
        appState.importConfiguration()
    }
    .keyboardShortcut("i", modifiers: [.command, .shift])
}

CommandMenu("Deployment") {
    Button("Deploy All Stacks") {
        Task { await appState.deployAll() }
    }
    .keyboardShortcut("d", modifiers: [.command, .shift])
    .disabled(appState.isDeploying)

    Button("Stop Deployment") {
        appState.stopDeployment()
    }
    .keyboardShortcut(".", modifiers: .command)
    .disabled(!appState.isDeploying)

    Divider()

    Button("View Logs") {
        appState.showLogs = true
    }
    .keyboardShortcut("l", modifiers: [.command, .option])

    Button("Run Health Check") {
        Task { await appState.runHealthCheck() }
    }
    .keyboardShortcut("h", modifiers: [.command, .shift])
}

CommandMenu("AWS") {
    Button("Switch Profile...") {
        appState.showProfileSwitcher = true
    }
    .keyboardShortcut("p", modifiers: [.command, .option])

    Button("Refresh Credentials") {
        Task { await appState.refreshCredentials() }
    }
    .keyboardShortcut("r", modifiers: [.command, .shift])
}

```

```

    }
  }
}

```

---

## UI Patterns (10 macOS Patterns)

The Swift Deployer follows these macOS design patterns:

1. **NavigationSplitView** - Sidebar + Content + Inspector
2. **Liquid Glass** - On navigation/controls only
3. **Toolbar-as-Command-Center** - Grouped actions + overflow
4. **Scroll Edge Effects** - Floating UI legibility
5. **Master List → Detail** - 3-level navigation
6. **Search as First-Class** - Toolbar trailing position
7. **Tables for Data** - Lists for collections
8. **Multi-Select + Context Menus** - Drag & drop support
9. **Full Menu Bar** - Keyboard shortcuts
10. **Settings Window + Inspectors** - macOS-native patterns # Admin Dashboard Reference

## Dashboard Architecture

### Technology Stack

- **Framework:** Next.js 14 (App Router)
- **Language:** TypeScript
- **Styling:** Tailwind CSS
- **Components:** shadcn/ui
- **Icons:** Lucide React
- **State:** React Query, Zustand
- **Forms:** React Hook Form, Zod

### Page Structure

```

apps/admin-dashboard/
  app/
    layout.tsx      # Root layout
    page.tsx        # Landing/login
    (dashboard)/
      layout.tsx    # Dashboard layout with sidebar
      page.tsx      # Overview dashboard
      [module]/page.tsx # Individual modules
  components/
    ui/             # shadcn/ui components
    workflow-editor/ # Visual workflow editor
    shared/         # Shared components
  lib/
    api.ts          # API client
    auth.ts         # Auth utilities

```

## Complete Page Inventory (43 Pages)

### Core Administration

Page	Route	Purpose
Overview	/	System health, key metrics, quick actions
Administrators	/administrators	Manage admin users, roles, permissions
Audit Logs	/audit-logs	View all system audit events
AWS Logs	/aws-logs	CloudWatch log viewer
Security	/security	Security settings, WAF, compliance
Settings	/settings	Platform configuration
System Config	/system-config	Advanced system settings

### AI & Models

Page	Route	Purpose
Models	/models	AI model configuration
Model Metadata	/model-metadata	Model capabilities & pricing
User Models	/user-models	Per-tenant model access
Providers	/providers	AI provider management

### Orchestration

Page	Route	Purpose
Orchestration	/orchestration	Workflow management
Orchestration Patterns	/orchestration-patterns	49 patterns library
Orchestration Editor	/orchestration-patterns/editor	Visual workflow editor

### Think Tank

Page	Route	Purpose
Think Tank	/thinktank	Session management
Cognition	/cognition	Cognitive settings
Cognitive Brain	/cognitive-brain	Brain configuration
Consciousness	/consciousness	Consciousness monitoring
Metacognition	/metacognition	Self-reflection settings
Planning	/planning	Goal planning
World Model	/world-model	World model state

## Billing & Cost

Page	Route	Purpose
Billing	<code>/billing</code>	Revenue, invoices, subscriptions
Cost	<code>/cost</code>	Cost analytics, budgets

## Analytics & Monitoring

Page	Route	Purpose
Analytics	<code>/analytics</code>	Usage analytics
Reports	<code>/reports</code>	Generated reports
Health	<code>/health</code>	System health dashboard
Deployments	<code>/deployments</code>	Deployment history

## AGI & Learning

Page	Route	Purpose
Agents	<code>/agents</code>	Autonomous agents
Learning	<code>/learning</code>	ML training data
ML Training	<code>/ml-training</code>	Model training jobs
Self-Improvement	<code>/self-improvement</code>	Self-improvement logs
Moral Compass	<code>/moral-compass</code>	Ethical guidelines
Feedback	<code>/feedback</code>	User feedback

## Collaboration & Features

Page	Route	Purpose
Time Machine	<code>/time-machine</code>	Historical state access
Storage	<code>/storage</code>	File storage management
Notifications	<code>/notifications</code>	Notification settings
Localization	<code>/localization</code>	i18n management
Configuration	<code>/configuration</code>	Dynamic configuration
Compliance	<code>/compliance</code>	Compliance dashboard
Geographic	<code>/geographic</code>	Geographic settings
Multi-Region	<code>/multi-region</code>	Multi-region config
Experiments	<code>/experiments</code>	A/B testing
Migrations	<code>/migrations</code>	Database migrations
Services	<code>/services</code>	Service status
Request Handler	<code>/request-handler</code>	Request routing

## Key Pages Detail

### Overview Dashboard (/)

**Metrics Displayed:** - Active tenants (24h) - Total API requests (24h) - Total tokens processed - Revenue (MTD) - Error rate - Average latency

**Quick Actions:** - View recent errors - Check provider health - Review pending approvals - Generate report

**Charts:** - Requests over time (7d) - Token usage by model - Revenue trend - Error rate trend

---

### Models Page (/models)

**Features:** - List all 106+ models - Filter by provider, capability - Enable/disable models - Set model pricing overrides - Configure fallback chains - View usage statistics

#### Model Card Display:

anthropic/claude-3-5-sonnet

Provider: Bedrock (primary), LiteLLM (fb)

Capabilities: reasoning, coding, vision

Context: 200K tokens

Pricing: \$3.00/\$15.00 per 1M tokens

Status: Enabled

Usage (24h): 1.2M tokens

[Configure] [Disable] [View Stats]

---

### Orchestration Patterns Page (/orchestration-patterns)

**Features:** - Browse 49 patterns by category - Search patterns - View pattern details - Edit pattern workflows - Create custom patterns - View execution statistics

**Pattern Categories Tabs:** - Consensus & Aggregation (7) - Debate & Deliberation (7) - Critique & Refinement (7) - Verification & Validation (7) - Decomposition (7) - Specialized Reasoning (7) - Multi-Model Routing (4) - Ensemble Methods (3)

#### Pattern Detail View:

AI Debate

[Edit] [Test]

Category: Debate & Deliberation

Quality Improvement: +25-40%

Typical Latency: High (10-30s)

Min Models: 3

#### Description:

Two AI models debate opposing positions while a third model judges the arguments and synthesizes a final answer.

#### Best For:

- Controversial topics
- Complex decisions
- Exploring multiple perspectives

#### Workflow Steps:

1. Generate Pro Argument (Claude)
2. Generate Con Argument (GPT-4o)
3. Judge Arguments (Claude - thinking mode)
4. Synthesize Final Answer

Executions (30d): 1,247 | Avg Quality: 0.89

---

### Visual Workflow Editor (/orchestration-patterns/editor)

**Features:** - Drag-and-drop workflow design - 16 method palette - Node connection editing - Step configuration (4 tabs) - Zoom, pan, fit controls - Test execution - Save/load workflows

**Method Palette (16 Methods):** | Method | Category | Description | |———|———|———|———|  
| Generate | Core | Generate text response | | Analyze | Core | Analyze input | | Transform | Core | Transform data | | Validate | Core | Validate output | | Critique | Refinement | Critique response | | Refine | Refinement | Improve response | | Decompose | Decomposition | Break into parts | | Synthesize | Aggregation | Combine results | | Judge | Evaluation | Evaluate quality | | Vote | Consensus | Majority voting | | Debate\_Pro | Debate | Pro argument | | Debate\_Con | Debate | Con argument | | Verify | Verification | Fact-check | | Search | External | Web search | | Execute\_Code | External | Run code | | Custom | Custom | Custom logic |

**Step Configuration Tabs:** 1. **General** - Name, order, model, output variable 2. **Parameters** - Method-specific parameters 3. **Advanced** - Conditions, iterations, dependencies 4. **Parallel** - Parallel execution settings, AGI selection

---

### Billing Page (/billing)

#### Sections:

**Revenue Overview:** - Monthly Recurring Revenue (MRR) - Annual Recurring Revenue (ARR)  
- Revenue growth % - Churn rate

**Subscription Management:** - Active subscriptions by tier - Upcoming renewals - Cancelled subscriptions - Trial conversions

**Credit Management:** - Total credits sold - Credits consumed - Credit purchase history - Low balance alerts

**Invoice Management:** - Generate invoices - View invoice history - Export to CSV - Send invoice reminders

---

## Analytics Page (/analytics)

### Dashboard Sections:

**Usage Analytics:** - Requests by model - Tokens by tenant - Peak usage times - Geographic distribution

**Performance Analytics:** - Latency percentiles (p50, p95, p99) - Error rates by endpoint - Provider availability - Cache hit rates

**Business Analytics:** - Cost per request - Revenue per tenant - Feature adoption - User engagement

**Custom Reports:** - Date range selection - Dimension grouping - Metric selection - Export options (CSV, PDF, JSON)

---

## Security Page (/security)

### Sections:

**Authentication:** - Cognito configuration - MFA enforcement - Session settings - Password policies

**API Security:** - Rate limiting rules - IP allowlists - API key management - Request validation

**Compliance:** - SOC2 status - HIPAA mode toggle - Data retention settings - Audit log retention

**WAF Configuration:** - Rule management - Blocked requests - Rate limit thresholds - Custom rules

---

## Think Tank Page (/thinktank)

**Features:** - View all sessions across tenants - Filter by domain, status, confidence - Session detail view - Step-by-step reasoning display - Cost and token tracking

### Session List View:

Session ID	Tenant	Domain	Steps	Conf	Cost
abc123...	Acme	Engineering	6	0.92	\$0.45
def456...	Beta	Research	8	0.87	\$0.72
ghi789...	Acme	Legal	4	0.95	\$0.28



## Session Detail View:

Problem: "Design a microservices architecture for 10M daily users"

Domain: Engineering | Complexity: High | Status: Completed

Step 1: Decompose [Claude 3.5] (conf: 0.94)

Identified 5 sub-problems

Duration: 2.3s | Tokens: 1,247

Step 2: Requirements Analysis [Claude + GPT-4o parallel] (conf: 0.91)

Synthesized from 2 models

Duration: 4.1s | Tokens: 3,892

Step 3-5: [...]

Step 6: Synthesize Final Solution [Claude 3.5 thinking] (conf: 0.89)

Generated comprehensive solution

Duration: 5.7s | Tokens: 2,156

Total: 6 steps | 18.2s | 12,453 tokens | \$0.45

Final Confidence: 0.89

---

## Component Library

### Shared Components

**DataTable:** - Sortable columns - Pagination - Row selection - Export functionality - Column visibility toggle

**StatusBadge:** - Status indicator with color - Configurable variants - Icon support

**MetricCard:** - Large number display - Trend indicator - Comparison to previous period

**Chart Components:** - LineChart (time series) - BarChart (comparisons) - PieChart (distributions) - AreaChart (cumulative)

**Form Components:** - Input with validation - Select with search - DateRangePicker - JSONEditor - CodeEditor

---

## API Integration

### API Client (lib/api.ts)

```
import { QueryClient } from '@tanstack/react-query';
```

```
const API_BASE = process.env.NEXT_PUBLIC_API_URL;
```

```
export const apiClient = {  
  // GET request
```

```

async get<T>(path: string): Promise<T> {
  const response = await fetch(`${API_BASE}${path}`, {
    headers: await getAuthHeaders(),
  });
  if (!response.ok) throw new APIError(response);
  return response.json();
},

// POST request
async post<T>(path: string, data: unknown): Promise<T> {
  const response = await fetch(`${API_BASE}${path}`, {
    method: 'POST',
    headers: {
      'Content-Type': 'application/json',
      ...(await getAuthHeaders()),
    },
    body: JSON.stringify(data),
  });
  if (!response.ok) throw new APIError(response);
  return response.json();
},

// React Query hooks
useModels: () => useQuery(['models'], () => apiClient.get('/admin/models')),
useTenants: () => useQuery(['tenants'], () => apiClient.get('/admin/tenants')),
useAnalytics: (range: string) =>
  useQuery(['analytics', range], () => apiClient.get(`/admin/analytics?range=${range}`)),
};

```

## Authentication (lib/auth.ts)

```

import { Amplify, Auth } from 'aws-amplify';

export async function getAuthHeaders(): Promise<Headers> {
  const session = await Auth.currentSession();
  return {
    'Authorization': `Bearer ${session.getIdToken().getJwtToken()}`,
  };
}

export function useAuth() {
  const [user, setUser] = useState<CognitoUser | null>(null);
  const [loading, setLoading] = useState(true);

  useEffect(() => {
    Auth.currentAuthenticatedUser()
      .then(setUser)
      .catch(() => setUser(null))
  });
}

```

```

        .finally(() => setLoading(false));
    }, []);

    return { user, loading, signIn, signOut };
}

```

## Compliance & Security Standards

### Overview

RADIANT implements comprehensive compliance frameworks to meet enterprise security requirements across multiple regulatory standards.

### Required Provider API Keys

RADIANT requires the following external AI provider API keys for deployment:

Provider	Secret Path	Purpose	Get Key
<b>Anthropic (Claude)</b>	radiant/providers/anthropic	Primary AI provider for Claude models	<a href="https://console.anthropic.com/settings/keys">https://console.anthropic.com/settings/keys</a>
<b>Groq</b>	radiant/providers/groq	Fast LPU inference for fallback	<a href="https://console.groq.com/keys">https://console.groq.com/keys</a>

### Deployment Flow

1. **Configure Keys in Deployer** - Enter API keys in the Swift Deployer “Required Provider API Keys” section
2. **Local Storage** - Keys are stored securely in macOS Keychain
3. **AWS Upload** - During deployment, keys are uploaded to AWS Secrets Manager
4. **Lambda Access** - Lambda functions retrieve keys from Secrets Manager at runtime

### Why These Providers Are Required

- **Anthropic (Claude)**: Primary provider for Claude 3.5 Sonnet, Claude Opus 4, and other Claude models via AWS Bedrock. Direct API access provides extended thinking and latest model features.
- **Groq**: Ultra-fast inference (100-200ms) on Llama and Mixtral models. Used as fallback when Bedrock is unavailable and for speed-critical applications.

## SOC 2 Type II Compliance

### Trust Service Criteria

Category	Controls	Implementation
<b>Security</b>	Access control, encryption, monitoring	Cognito, KMS, CloudWatch
<b>Availability</b>	Redundancy, failover, SLAs	Multi-AZ, auto-scaling
<b>Processing Integrity</b>	Data validation, error handling	Input validation, checksums
<b>Confidentiality</b>	Data classification, encryption	RLS, AES-256, TLS 1.3
<b>Privacy</b>	Data handling, consent	GDPR controls, retention policies

## Key Controls

### 1. Access Management

- Multi-factor authentication (MFA) required for admins
- Role-based access control (RBAC)
- API key rotation policies
- Session timeout enforcement

### 2. Encryption

- At rest: AES-256 via AWS KMS
- In transit: TLS 1.3 minimum
- Database: Aurora encryption enabled
- Secrets: AWS Secrets Manager

### 3. Audit Logging

- All API requests logged
- Admin actions tracked
- CloudTrail for AWS operations
- 90-day retention minimum

## HIPAA Compliance

### Protected Health Information (PHI) Handling

RADIANT supports HIPAA-compliant deployments with enhanced controls:

Requirement	Implementation
<b>Access Controls</b>	User authentication, authorization, audit
<b>Audit Controls</b>	Complete activity logging, tamper-evident
<b>Integrity Controls</b>	Data validation, checksums, versioning
<b>Transmission Security</b>	TLS 1.3, encrypted channels only

## HIPAA Mode Features

When HIPAA mode is enabled:

```

interface HIPAAConfig {
  enabled: boolean;
  phiDetection: boolean;      // Scan for PHI in requests
  enhancedLogging: boolean;   // Additional audit details
  dataRetentionDays: number;  // Configurable retention
  encryptionRequired: boolean; // Force encryption
  accessReviewDays: number;   // Periodic access review
}

```

## PHI Sanitization

```

// Automatic PHI detection and handling
export class PHISanitizationService {
  private patterns = [
    /\b\d{3}-\d{2}-\d{4}\b/,    // SSN
    /\b\d{9}\b/,               // MRN
    /\b[A-Z]{2}\d{6,8}\b/,     // License numbers
    // ... additional patterns
  ];

  async sanitize(input: string): Promise<SanitizedInput> {
    // Detect and redact PHI before processing
    let sanitized = input;
    for (const pattern of this.patterns) {
      sanitized = sanitized.replace(pattern, '[REDACTED]');
    }
    return { original: input, sanitized, phiDetected: sanitized !== input };
  }
}

```

---

## GDPR Compliance

### Data Subject Rights

RADIANT implements all required GDPR data subject rights:

Right	Implementation	API Endpoint
<b>Right to Access</b>	Export all user data	GET /api/gdpr/export
<b>Right to Rectification</b>	Update personal data	PATCH /api/users/{id}
<b>Right to Erasure</b>	Delete all user data	DELETE /api/gdpr/erase
<b>Right to Portability</b>	Export in machine-readable format	GET /api/gdpr/export?format=json
<b>Right to Object</b>	Opt-out of processing	POST /api/gdpr/object

Right	Implementation	API Endpoint
<b>Right to Restrict</b>	Limit processing	POST /api/gdpr/restrict

## Data Processing

```
interface GDPRDataRequest {
  subjectId: string;           // User identifier
  requestType: 'access' | 'rectification' | 'erasure' | 'portability' | 'object' | 'restrict';
  requestedBy: string;         // Requester (user or DPO)
  verificationMethod: string;  // How identity was verified
  deadline: Date;              // 30-day compliance deadline
}

export class GDPRService {
  async handleDataRequest(request: GDPRDataRequest): Promise<GDPRResponse> {
    // Log the request
    await this.auditLogger.log('gdpr_request', request);

    switch (request.requestType) {
      case 'access':
        return this.exportUserData(request.subjectId);
      case 'erasure':
        return this.eraseUserData(request.subjectId);
      case 'portability':
        return this.exportPortableData(request.subjectId);
      // ... other handlers
    }
  }

  async eraseUserData(userId: string): Promise<void> {
    // Cascade delete across all tables
    await this.db.transaction(async (tx) => {
      await tx.delete('thinktank_steps').where('session_id', 'in',
        tx.select('id').from('thinktank_sessions').where('user_id', userId));
      await tx.delete('thinktank_sessions').where('user_id', userId);
      await tx.delete('usage_records').where('user_id', userId);
      await tx.delete('api_keys').where('user_id', userId);
      await tx.delete('users').where('id', userId);
    });
  }
}
```

## Consent Management

```
interface ConsentRecord {
  userId: string;
```

```

    consentType: 'marketing' | 'analytics' | 'ai_training' | 'data_sharing';
    granted: boolean;
    grantedAt: Date;
    ipAddress: string;
    userAgent: string;
    version: string; // Consent policy version
}

// All processing requires valid consent
async function checkConsent(userId: string, purpose: string): Promise<boolean> {
    const consent = await db.query(
        'SELECT granted FROM consent_records WHERE user_id = $1 AND consent_type = $2',
        [userId, purpose]
    );
    return consent?.granted === true;
}

```

## Data Retention

Data Type	Retention Period	Basis
User accounts	Until deletion requested	Contract
Session data	90 days	Legitimate interest
Audit logs	7 years	Legal requirement
Usage analytics	2 years	Legitimate interest
AI training data	Until consent withdrawn	Consent

## ISO 27001 Compliance

### Information Security Management System (ISMS)

RADIANT's infrastructure aligns with ISO 27001:2022 requirements:

#### Annex A Controls

##### A.5 Organizational Controls

Control	Description	Implementation
A.5.1	Policies for information security	Documented security policies
A.5.2	Information security roles	Defined RACI matrix
A.5.3	Segregation of duties	Role-based access, dual approval
A.5.7	Threat intelligence	AWS GuardDuty, threat feeds
A.5.15	Access control	Cognito + IAM + RLS
A.5.23	Information security for cloud	AWS Well-Architected
A.5.29	Information security during disruption	DR procedures

## A.6 People Controls

Control	Description	Implementation
A.6.1	Screening	Background checks for admins
A.6.3	Information security awareness	Training programs
A.6.5	Responsibilities after termination	Access revocation procedures

## A.7 Physical Controls

Control	Description	Implementation
A.7.1	Physical security perimeters	AWS data center security
A.7.4	Physical security monitoring	AWS compliance certifications

## A.8 Technological Controls

Control	Description	Implementation
A.8.1	User endpoint devices	MDM for admin devices
A.8.2	Privileged access rights	IAM policies, MFA required
A.8.3	Information access restriction	RLS, tenant isolation
A.8.4	Access to source code	GitHub branch protection
A.8.5	Secure authentication	Cognito, JWT, API keys
A.8.7	Protection against malware	WAF, input validation
A.8.9	Configuration management	CDK, Infrastructure as Code
A.8.10	Information deletion	GDPR erasure, retention policies
A.8.11	Data masking	PHI sanitization, PII redaction
A.8.12	Data leakage prevention	DLP policies, egress controls
A.8.15	Logging	CloudWatch, audit trails
A.8.16	Monitoring activities	CloudWatch alarms, dashboards
A.8.20	Networks security	VPC, security groups, NACLs
A.8.22	Segregation of networks	Private subnets, VPC endpoints
A.8.24	Use of cryptography	KMS, TLS 1.3, AES-256
A.8.25	Secure development lifecycle	Code review, security scanning
A.8.28	Secure coding	OWASP guidelines, linting

## Risk Assessment Matrix

Risk Category	Likelihood	Impact	Controls
Data breach	Low	Critical	Encryption, RLS, monitoring
Service outage	Medium	High	Multi-AZ, auto-scaling, DR
Unauthorized access	Low	Critical	MFA, RBAC, audit logging
Insider threat	Low	High	Segregation, dual approval
Supply chain attack	Low	High	Dependency scanning, SBOMs



## Security Architecture

### Defense in Depth

INTERNET

AWS WAF

Rate limiting, SQL injection, XSS protection

CloudFront / ALB

TLS 1.3 termination

API Gateway

JWT validation, API key auth

VPC (Private)

Lambda Functions

Input validation, RLS context

Aurora PostgreSQL

Row-Level Security, Encryption

### Two-Person Approval

Sensitive operations require dual admin approval:

```
interface ApprovalRequest {
    id: string;
    requesterId: string;
    actionType: 'delete_tenant' | 'modify_billing' | 'grant_super_admin' |
                'bulk_export' | 'disable_security';
    resourceId: string;
    payload: Record<string, unknown>;
    status: 'pending' | 'approved' | 'rejected' | 'expired';
    requiredApprovals: number; // Usually 2
```

```

    approvals: Approval[];
    expiresAt: Date;
}

// Cannot approve own requests
async function approveRequest(requestId: string, approverId: string) {
    const request = await getRequest(requestId);

    if (request.requesterId === approverId) {
        throw new Error('Cannot approve own request');
    }

    if (request.approvals.some(a => a.approverId === approverId)) {
        throw new Error('Already approved');
    }

    // Add approval
    request.approvals.push({ approverId, approvedAt: new Date() });

    // Execute if threshold met
    if (request.approvals.length >= request.requiredApprovals) {
        await executeApprovedAction(request);
    }
}

```

---

## Audit Logging

### Log Structure

```

interface AuditLog {
    id: string;
    timestamp: Date;
    tenantId: string;
    userId: string;
    adminId?: string;
    action: string;
    resourceType: string;
    resourceId: string;
    ipAddress: string;
    userAgent: string;
    requestId: string;
    oldValue?: Record<string, unknown>;
    newValue?: Record<string, unknown>;
    result: 'success' | 'failure';
    errorMessage?: string;
}

```

## Logged Actions

- All authentication events (login, logout, MFA)
- All API requests with parameters
- All data modifications (create, update, delete)
- All admin actions
- All security events (failed auth, rate limits)
- All GDPR requests
- All compliance-related operations

## Log Retention

Log Type	Retention	Storage
API Access	90 days	CloudWatch
Security Events	1 year	S3 + Glacier
Audit Trail	7 years	S3 + Glacier
GDPR Requests	7 years	Aurora