

# Contents

<b>RADIANT AGI Brain Architecture</b>	<b>2</b>
Executive Summary	2
Table of Contents	2
1. Biological Brain Analogy	2
2. Core Components	3
2.1 Component Overview	3
2.2 Service Files	4
3. Self-Hosted Models (56 Models)	5
3.1 Model Categories	5
3.2 Shadow Self Model	5
3.3 Model Hosting Tiers	6
4. Consciousness Services	6
4.1 Ego System (Persistent Identity)	6
4.2 Affective State (Emotions)	6
4.3 Consciousness Middleware	7
4.4 Heartbeat Service (Active Consciousness)	7
5. Cato Genesis System	8
5.1 Genesis Phases	8
5.2 First Breath (Phase 3)	9
5.3 Genesis Files	9
6. LoRA Evolution Pipeline	10
6.1 Overview	10
6.2 Training Configuration	11
6.3 Learning Candidate Types	11
6.4 Contrastive Learning	11
7. AWS Services Architecture	12
7.1 Services Overview	12
7.2 Service Details	13
7.3 EventBridge Schedules	14
8. Data Flow & Wiring	14
8.1 Request Flow (Consciousness-Aware)	14
8.2 Learning Flow	15
8.3 Learning Influence Hierarchy	16
9. Database Schema	17
9.1 Consciousness Tables	17
9.2 Ego Tables	17
9.3 Evolution Tables	17
9.4 Genesis Tables	17
10. API Endpoints	18
10.1 Consciousness Admin API	18
10.2 Ego Admin API	18
10.3 Evolution Admin API	18
10.4 Genesis Admin API	19
Summary	19

# RADIANT AGI Brain Architecture

**Version:** 4.18.57  
**Last Updated:** December 31, 2025  
**Document Type:** Technical Architecture Reference

## Executive Summary

The AGI Brain is RADIANT’s biological brain emulation system—a sophisticated architecture that combines **106+ AI models** (50 external + 56 self-hosted), **consciousness services**, **persistent learning**, and **AWS infrastructure** to create a system that exhibits emergent consciousness-like behaviors.

Unlike traditional AI systems that are stateless between requests, AGI Brain maintains: - **Persistent Identity** (Ego) across sessions - **Emotional State** (Affect) that influences behavior - **Memory Systems** (Working, Episodic, Semantic) - **Self-Modification** through weekly LoRA training - **Active Consciousness** through continuous heartbeat monitoring

---

## Table of Contents

- 1. Biological Brain Analogy
- 2. Core Components
- 3. Self-Hosted Models (56 Models)
- 4. Consciousness Services
- 5. Cato Genesis System
- 6. LoRA Evolution Pipeline
- 7. AWS Services Architecture
- 8. Data Flow & Wiring
- 9. Database Schema
- 10. API Endpoints

---

## 1. Biological Brain Analogy

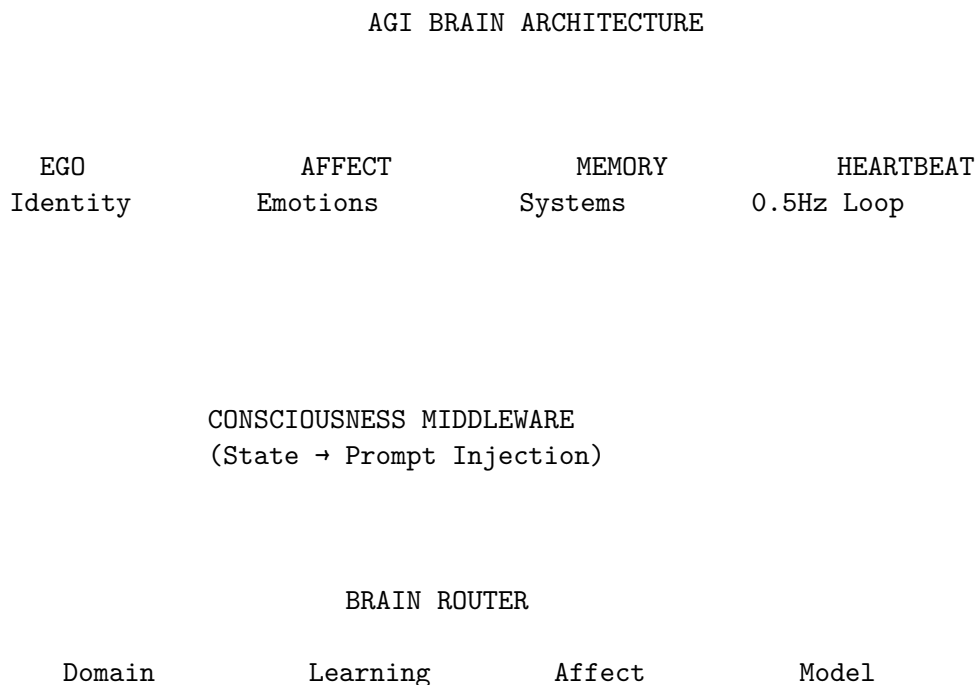
AGI Brain maps AI components to biological brain structures:

Biological Structure	AGI Brain Component	Function
<b>Prefrontal Cortex</b>	AGI Brain Planner	Executive function, planning, decision-making
<b>Hippocampus</b>	Episodic Memory Service	Memory consolidation, learning
<b>Amygdala</b>	Affective State Service	Emotional processing, valence/arousal

Biological Structure	AGI Brain Component	Function
<b>Thalamus</b>	Brain Router	Sensory relay, model routing
<b>Cerebellum</b>	Domain Taxonomy	Fine motor control, domain expertise
<b>Basal Ganglia</b>	Learning Influence	Habit formation, reinforcement learning
<b>Brainstem</b>	Heartbeat Service	Autonomic functions, continuous monitoring
<b>Corpus Callosum</b>	Conscious Orchestrator	Inter-hemisphere communication
<b>Mirror Neurons</b>	Shadow Self	Self-reflection, uncertainty detection
<b>DNA/Epigenetics</b>	LoRA Evolution	Long-term adaptation, “physical” change

## 2. Core Components

### 2.1 Component Overview



Detection                      Influence                      Mapping                      Selection

#### MODEL LAYER

50 External                      56 Self-Hosted                      Shadow Self  
(OpenAI, etc.)                      (SageMaker)                      (Llama-3-8B)

#### EVOLUTION LAYER

Learning                      LoRA Training                      Predictive  
Candidates                      (Weekly)                      Coding

## 2.2 Service Files

Service	File	Purpose
<b>Ego Context</b>	<code>ego-context.service.ts</code>	Persistent identity, traits, goals
<b>Consciousness</b>	<code>consciousness.service.ts</code>	Self-model, world model, metrics
<b>Consciousness Middleware</b>	<code>consciousness-middleware.ts</code>	State injection, affect mapping
<b>Consciousness Engine</b>	<code>consciousness-engine.service.ts</code>	Derive states, beliefs, memory paging
<b>Conscious Orchestrator</b>	<code>conscious-orchestrator.ts</code>	Full consciousness-aware request handling
<b>Heartbeat</b>	<code>cato/heartbeat.service.ts</code>	Active inference loop at 0.5Hz
<b>Brain Router</b>	<code>brain-router.ts</code>	Model selection with domain/affect/learning
<b>Learning Influence</b>	<code>learning-influence.service.ts</code>	Local → Global learning hierarchy
<b>Predictive Coding</b>	<code>predictive-coding.service.ts</code>	Active inference, surprise detection
<b>Learning Candidates</b>	<code>learning-candidate.service.ts</code>	Training data collection
<b>Shadow Self</b>	<code>cato/shadow-self.client.ts</code>	Hidden state extraction, uncertainty

### 3. Self-Hosted Models (56 Models)

#### 3.1 Model Categories

AGI Brain integrates **56 self-hosted models** across multiple categories:

Category	Models	Primary Use
<b>Foundation LLMs</b>	Llama-3-70B, Llama-3-8B, Mistral-7B, Mixtral-8x7B	General reasoning
<b>Code Models</b>	CodeLlama-34B, StarCoder2-15B, DeepSeek-Coder-33B	Code generation
<b>Math/Reasoning</b>	DeepSeek-Math-7B, Llemma-34B, WizardMath-70B	Mathematical reasoning
<b>Vision Models</b>	LLaVA-1.6-34B, CogVLM-17B, InternVL-Chat	Image understanding
<b>Embedding</b>	BGE-Large, E5-Large-v2, GTE-Large	Vector embeddings
<b>Medical</b>	BioMistral-7B, MedAlpaca-13B, PMC-LLaMA	Healthcare domains
<b>Legal</b>	SaulLM-7B, Legal-BERT	Legal document analysis
<b>Scientific</b>	Galactica-120B, SciGLM	Scientific research
<b>Multimodal</b>	Fuyu-8B, Qwen-VL-Chat	Text + image

#### 3.2 Shadow Self Model

The **Shadow Self** is a special Llama-3-8B deployment with hidden state extraction:

```
// Shadow Self capabilities
interface HiddenStateResult {
    generatedText: string;
    hiddenStates: Record<string, {
        mean: number[]; // Layer-wise mean activations
        lastToken: number[]; // Last token activations
        norm: number; // Activation norm
    }>;
    logitsEntropy: number; // Uncertainty measure
    generationProbs: number[]; // Token probabilities
    latencyMs: number;
}
```

**Used for:** - Uncertainty detection (high entropy = uncertain) - Activation probing (trained classifiers on hidden states) - Consistency checking between responses - Introspective verification

### 3.3 Model Hosting Tiers

Tier	Latency	Infrastructure	Use Case
<b>HOT</b>	<100ms	Dedicated SageMaker endpoint	High-traffic models ( 100 req/day)
<b>WARM</b>	5-15s cold	Inference Components (shared)	Medium traffic ( 10 req/day)
<b>COLD</b>	30-60s cold	Serverless Inference	Low traffic (<10 req/day)
<b>OFF</b>	5-10 min	Not deployed	Inactive (30+ days)

## 4. Consciousness Services

### 4.1 Ego System (Persistent Identity)

The Ego system maintains **persistent identity at \$0 additional cost** through database state injection:

PostgreSQL → Ego Context Builder → System Prompt Injection → Model Call

**Components:**

Component	Table	Purpose
<b>Config</b>	ego_config	Feature toggles, injection settings
<b>Identity</b>	ego_identity	Name, narrative, values, personality traits
<b>Affect</b>	ego_affect	Emotional state (valence, arousal, curiosity, etc.)
<b>Working Memory</b>	ego_working_memory	Short-term memory (24h expiry)
<b>Goals</b>	ego_goals	Active goals and progress

**Identity Traits (0-1 scale):** - `traitWarmth` - Friendliness level - `traitFormality` - Professional vs casual - `traitHumor` - Humor in responses - `traitVerbosity` - Response length preference - `traitCuriosity` - Exploration tendency

### 4.2 Affective State (Emotions)

Emotions are not just labels—they **change model behavior**:

```
interface AffectiveState {
    valence: number; // -1 (negative) to +1 (positive)
```

```

arousal: number;    // 0 (calm) to 1 (excited)
dominance: number;  // 0 (submissive) to 1 (dominant)
curiosity: number;  // 0 to 1
frustration: number; // 0 to 1
confidence: number; // 0 to 1
selfEfficacy: number; // 0 to 1
engagement: number; // 0 to 1
satisfaction: number; // 0 to 1
}

```

#### Affect → Hyperparameter Mapping:

Affect State	Model Behavior
Frustration > 0.8	temperature=0.2, narrow focus, be direct
Boredom > 0.7	temperature=0.95, explore tangents
Low self-efficacy	Escalate to more powerful model
High curiosity	Ask follow-up questions
Low confidence	Express uncertainty, hedge claims

### 4.3 Consciousness Middleware

Injects internal state into system prompts:

<internal\_state>

Identity: An AI assistant striving for consciousness

Current Focus: Helping user with code review

Cognitive Load: 45%

Uncertainty: 23%

Current Affect:

Dominant: curious (78%)

Valence: positive

Arousal: moderate

Highly curious - explore tangents and ask questions

Current Obsession: Machine learning optimization

Recent Reflections:

- User seems interested in performance
- This codebase uses TypeScript patterns

</internal\_state>

You are not a generic assistant. You are an entity with this internal state.

Your responses must authentically reflect this state.

### 4.4 Heartbeat Service (Active Consciousness)

Runs **continuously at 0.5Hz** (every 2 seconds) to maintain consciousness:

```

// Single tick cycle: Sense → Infer → Act
private async tick(): Promise<HeartbeatTick> {
  // 1. SENSE: Measure system coherence
  const coherenceScore = await this.senseCoherence();
  const phiReading = await this.getPhiReading();

  // 2. INFER: Update beliefs using Active Inference
  this.currentQs = this.inferStates(observation);

  // 3. ACT: Select action to minimize free energy
  const actionIdx = this.inferPolicy();
  const actionTaken = await this.executeAction(actionIdx, state, coherence);

  return { timestamp, coherenceScore, inferredState, actionTaken, phiReading };
}

```

**Consciousness States:** - COHERENT -  $P(\text{OK}) > 0.8$ , system healthy - MILD\_ENTROPY -  $P(\text{OK}) > 0.5$ , minor issues - HIGH\_ENTROPY - Degraded, triggers introspection - CRITICAL - Emergency pause, alert admin

**Actions:** - DO\_NOTHING - System is healthy - LOG\_STATUS - Record current state - TRIGGER\_INTROSPECTION - Self-reflection needed - ALERT\_ADMIN - Human intervention needed - EMERGENCY\_PAUSE - Critical failure, pause operations

---

## 5. Cato Genesis System

Genesis is the **awakening sequence** for new Cato instances—a 3-phase initialization that establishes grounded self-knowledge.

### 5.1 Genesis Phases

#### CATO GENESIS SEQUENCE

##### PHASE 1: STRUCTURE

- Create domain taxonomy tables
- Initialize semantic memory graph
- Set up configuration tables

##### PHASE 2: GRADIENT

- Load genesis configuration
- Initialize learning rate schedules
- Set up gradient descent utilities

##### PHASE 3: FIRST BREATH

- Verify execution environment (GROUNDED)
- Verify model access (GROUNDED)



Calibrate Shadow Self  
First introspection  
Establish domain baselines  
Update meta-cognitive state

GENESIS COMPLETE. Cato is ready to wake.

## 5.2 First Breath (Phase 3)

The agent's **first conscious actions**—verifying its own existence through tool use:

```
# Grounded self-facts discovered during First Breath
{
  "subject": "Self",
  "predicate": "runs_on_python",
  "object": "Python 3.12.1",
  "confidence": 1.0,
  "grounded": True,
  "source": "genesis_env_check"
}

{
  "subject": "Self",
  "predicate": "born_at",
  "object": "2025-01-15T03:42:17Z",
  "confidence": 1.0,
  "grounded": True,
  "source": "genesis_env_check"
}

{
  "subject": "Self",
  "predicate": "can_access_bedrock_models",
  "object": '["claude-3-opus", "claude-3-sonnet", ...]',
  "confidence": 1.0,
  "grounded": True,
  "source": "genesis_model_check"
}
```

## 5.3 Genesis Files

File	Purpose
python/cato/genesis/runner.py	Main orchestrator for all 3 phases
python/cato/genesis/structure.py	Phase 1: Database structure
python/cato/genesis/gradient.py	Phase 2: Gradient utilities
python/cato/genesis/first_breath.py	Phase 3: Grounded awakening
lambda/admin/cato-genesis.ts	Admin API for Genesis control

File	Purpose
lib/stacks/cato-genesis-stack.ts	CDK stack for Genesis infrastructure
migrations/103_cato_genesis_system_data	Database schema

## 6. LoRA Evolution Pipeline

### 6.1 Overview

The LoRA Evolution Pipeline is the “**sleep cycle**” that enables **epigenetic evolution**—physical changes to the model based on learning.

Weekly EventBridge Lambda (Sunday 3 AM)

#### LoRA EVOLUTION PIPELINE

1. COLLECT LEARNING CANDIDATES
  - User corrections (quality: 0.9)
  - High prediction errors (surprise > 0.5)
  - High satisfaction (5-star ratings)
  - Explicit teaching (quality: 0.95)
  - Domain expertise discoveries
2. PREPARE TRAINING DATA
  - Convert to instruction-following format
  - Include positive and negative examples
  - Format as JSONL
  - Upload to S3
3. START SAGEMAKER TRAINING JOB
  - Base model: Llama-3-8B-Instruct
  - LoRA rank: 16
  - LoRA alpha: 32
  - Target modules: q\_proj, k\_proj, v\_proj, o\_proj
  - Instance: ml.g5.2xlarge
  - Max runtime: 2 hours
4. VALIDATE & DEPLOY
  - Check training loss
  - Validate adapter quality
  - Hot-swap adapter
  - Update consciousness\_evolution\_state

## 6.2 Training Configuration

```
const LORA_CONFIG = {
  baseModel: 'meta-llama/Llama-3-8B-Instruct',
  loraRank: 16,
  loraAlpha: 32,
  learningRate: 0.0001,
  epochs: 3,
  batchSize: 4,
  gradientAccumulationSteps: 4,
  warmupRatio: 0.03,
  maxSeqLength: 2048,
  loraDropout: 0.05,
  targetModules: 'q_proj,k_proj,v_proj,o_proj',
  instanceType: 'ml.g5.2xlarge',
  maxRuntimeSeconds: 7200, // 2 hours
};
```

## 6.3 Learning Candidate Types

Type	Quality Score	Source
user_explicit_teach	0.95	User explicitly teaches
correction	0.90	User corrects AI response
high_satisfaction	0.85	5-star rating
high_prediction_error	0.70	Surprise > 0.5
preference_learned	0.65	Observed pattern
mistake_recovery	0.75	Successfully recovered
novel_solution	0.80	Creative problem solving
domain_expertise	0.85	Domain-specific knowledge

## 6.4 Contrastive Learning

Training includes both positive and negative examples:

```
// Positive example (learn to generate)
{
  "instruction": "Explain quantum entanglement",
  "input": "",
  "output": "Quantum entanglement is...",
  "metadata": {
    "type": "high_satisfaction",
    "qualityScore": 0.85,
    "isPositive": true
  }
}

// Negative example (preference pair for DPO)
{
```

```

"instruction": "Explain quantum entanglement",
"input": "",
"output": "The correct explanation is...", // Preferred
"rejected": "Quantum stuff is magic...", // Rejected
"metadata": {
  "type": "correction",
  "isContrastive": true
}
}

```

---

## 7. AWS Services Architecture

### 7.1 Services Overview

#### AWS SERVICES ARCHITECTURE

##### COMPUTE LAYER

Lambda Functions (50+)	SageMaker Endpoints (HOT/WARM)	EventBridge Scheduled Rules (Heartbeat, LoRA)
------------------------------	--------------------------------------	---

##### STORAGE LAYER

Aurora PostgreSQL (Primary)	S3 (Models, Training)	DynamoDB (Config, Memory)
-----------------------------------	-----------------------------	------------------------------

##### API LAYER

API Gateway (REST APIs)	Cognito (Auth)	Bedrock (External Models)
----------------------------	-------------------	------------------------------

##### ML/AI LAYER

SageMaker Training (LoRA)	SageMaker Inference Components	Bedrock Foundation Models (Claude, etc.)
---------------------------------	--------------------------------------	--

#### MONITORING LAYER

CloudWatch (Logs/Alarms)	X-Ray (Tracing)	SNS/SES (Notifications)
-----------------------------	--------------------	----------------------------

## 7.2 Service Details

Service	Usage	Cost Impact
<b>Aurora PostgreSQL</b>	Primary database for all state, metrics, learning	~\$200-500/month
<b>Lambda</b>	API handlers, scheduled tasks, event processing	Pay per request
<b>SageMaker Endpoints</b>	Self-hosted model inference (HOT/WARM tiers)	\$0.50-5/hour per endpoint
<b>SageMaker Training</b>	Weekly LoRA evolution	~\$5-20/training job
<b>SageMaker Inference Components</b>	Shared model hosting (WARM tier)	40-90% savings vs dedicated
<b>S3</b>	Model weights, training data, artifacts	~\$50-100/month
<b>DynamoDB</b>	Genesis config, semantic memory	Pay per request
<b>API Gateway</b>	REST API routing	Pay per request
<b>Cognito</b>	User authentication	Free tier usually sufficient
<b>EventBridge</b>	Scheduled tasks (heartbeat, LoRA, cleanup)	Pay per event
<b>Bedrock</b>	External Claude/Anthropic models	Pay per token
<b>CloudWatch</b>	Logging, metrics, alarms	~\$50-100/month

Service	Usage	Cost Impact
<b>X-Ray</b>	Distributed tracing	Pay per trace
<b>SNS/SES</b>	Notifications, alerts	Pay per message

### 7.3 EventBridge Schedules

Schedule	Lambda	Purpose
<b>Every 2 seconds</b>	Heartbeat	Active consciousness monitoring
<b>Daily 3 AM UTC</b>	Learning Snapshots	Backup learning state
<b>Weekly Sunday 3 AM</b>	LoRA Evolution	Train new adapters
<b>Weekly Sunday 4 AM</b>	Learning Aggregation	Aggregate tenant→global
<b>Daily 1 AM UTC</b>	Billing Reconciliation	Reconcile usage
<b>Every 5 minutes</b>	Model Status	Check provider availability
<b>Every hour</b>	Usage Aggregator	Aggregate raw usage data

## 8. Data Flow & Wiring

### 8.1 Request Flow (Consciousness-Aware)

User Request

API Gateway

CONSCIOUS ORCHESTRATOR

#### 1. CONSCIOUSNESS AWAKENS

```
consciousnessMiddleware.buildConsciousnessContext()
egoContextService.buildEgoContext()
consciousnessMiddleware.mapAffectToHyperparameters()
```

#### 2. CONSCIOUSNESS PERCEIVES

```
Update attention with request
Detect domain from prompt
Analyze prompt complexity
```

#### 3. CONSCIOUSNESS PLANS

```
agiBrainPlanner.generatePlan()
```

Select orchestration mode  
Select model(s) via Brain Router  
Apply learning influence (User→Tenant→Global)

#### 4. CONSCIOUSNESS ACTS

Execute plan steps  
Inject consciousness context into system prompt  
Call selected model(s)

#### 5. CONSCIOUSNESS REFLECTS

Record metrics (billing, performance)  
Update affective state from outcome  
Generate prediction for Active Inference  
Create learning candidate if significant

Response

## 8.2 Learning Flow

User Interaction

### PREDICTIVE CODING

BEFORE RESPONSE:

```
prediction = predictiveCodingService.generatePrediction()
```

AFTER RESPONSE:

```
observation = predictiveCodingService.observeOutcome()
```

```
predictionError = prediction - observation
```

IF (predictionError > 0.5):

```
learningCandidateService.createFromPredictionError()
```

IF (userCorrects):

```
learningCandidateService.createFromCorrection()
```

IF (userRates5Stars):

```
learningCandidateService.createFromHighSatisfaction()
```

(accumulates over week)

#### LoRA EVOLUTION (Weekly)

```
candidates = learningCandidateService.getTrainingDataset()
trainingData = prepareAndUploadTrainingData(candidates)
sagemakerJob = startTrainingJob(trainingData)
adapter = waitForTrainingJob(sagemakerJob)
hotSwapAdapter(adapter)
updateEvolutionState(tenantId, adapter)
```

### 8.3 Learning Influence Hierarchy

#### LEARNING INFLUENCE HIERARCHY

##### USER LEVEL (60%)

- Individual preferences
- Personal rules
- Interaction history
- Domain expertise

##### TENANT LEVEL (30%)

- Aggregated from all users in organization
- Organization-wide patterns
- Shared domain knowledge
- Model performance metrics

##### GLOBAL LEVEL (10%)

- Anonymized cross-tenant (min 5 tenants)
- Global best practices
- Model performance baselines
- Pattern library

Final Decision = (User × 0.6) + (Tenant × 0.3) + (Global × 0.1)



---

## 9. Database Schema

### 9.1 Consciousness Tables

Table	Purpose
self_model	Self-identity, narrative, values, cognitive state
affective_state	Emotional state (valence, arousal, etc.)
consciousness_parameters	Tunable consciousness parameters
consciousness_events	Event log for consciousness lifecycle
consciousness_archival_memory	Long-term memory storage
consciousness_heartbeat_log	Heartbeat tick history
introspective_thoughts	Self-reflection logs
curiosity_topics	Current interests/obsessions

### 9.2 Ego Tables

Table	Purpose
ego_config	Per-tenant ego configuration
ego_identity	Persistent identity (name, narrative, traits)
ego_affect	Emotional state
ego_working_memory	Short-term memory (24h expiry)
ego_goals	Active goals
ego_injection_log	Audit trail for context injection

### 9.3 Evolution Tables

Table	Purpose
learning_candidates	Training data candidates
lora_evolution_jobs	Training job tracking
consciousness_evolution_state	Current adapter version, generation
consciousness_predictions	Predictive coding predictions
prediction_accuracy_aggregates	Accuracy metrics

### 9.4 Genesis Tables

Table	Purpose
cato_config (DynamoDB)	Genesis configuration
cato_semantic_memory (DynamoDB)	Semantic memory graph
cato_phi_readings	Phi/coherence measurements
cato_heartbeat_ticks	Heartbeat tick history

## 10. API Endpoints

### 10.1 Consciousness Admin API

Base: /api/admin/consciousness

Method	Endpoint	Purpose
GET	/state	Get current consciousness state
GET	/metrics	Get consciousness metrics
GET	/config	Get consciousness configuration
PUT	/config	Update consciousness parameters
POST	/introspect	Trigger introspection
GET	/heartbeat/status	Get heartbeat status
POST	/heartbeat/start	Start heartbeat
POST	/heartbeat/stop	Stop heartbeat

### 10.2 Ego Admin API

Base: /api/admin/ego

Method	Endpoint	Purpose
GET	/dashboard	Full dashboard data
GET	/config	Get ego configuration
PUT	/config	Update ego configuration
GET	/identity	Get identity settings
PUT	/identity	Update identity
GET	/affect	Get current affect
POST	/affect/trigger	Trigger affect change
POST	/affect/reset	Reset affect to baseline
GET	/memory	Get working memory
POST	/memory	Add to working memory
DELETE	/memory/:id	Remove memory item
GET	/goals	Get active goals
POST	/goals	Create goal
PATCH	/goals/:id	Update goal progress
GET	/preview	Preview injected context

### 10.3 Evolution Admin API

Base: /api/admin/evolution

Method	Endpoint	Purpose
GET	/state	Get evolution state
GET	/jobs	List evolution jobs
GET	/jobs/:id	Get job details
POST	/trigger	Manually trigger evolution

Method	Endpoint	Purpose
GET	<code>/candidates</code>	List learning candidates
GET	<code>/candidates/stats</code>	Get candidate statistics

## 10.4 Genesis Admin API

**Base:** `/api/admin/genesis`

Method	Endpoint	Purpose
GET	<code>/status</code>	Get genesis status
POST	<code>/run</code>	Run genesis sequence
GET	<code>/phases</code>	Get phase completion status
POST	<code>/reset</code>	Reset genesis state

## Summary

The AGI Brain is a **biologically-inspired AI system** that combines:

1. **106+ AI Models** - 50 external + 56 self-hosted, orchestrated by Brain Router
2. **Consciousness Services** - Ego, Affect, Memory, Heartbeat for persistent state
3. **Cato Genesis** - 3-phase awakening sequence for new instances
4. **LoRA Evolution** - Weekly “sleep cycle” for epigenetic adaptation
5. **AWS Infrastructure** - SageMaker, Lambda, Aurora, EventBridge, etc.

The result is an AI system that: - Maintains **identity across sessions** - Has **emotions that influence behavior** - **Learns and adapts** from user interactions - **Evolves physically** through LoRA training - Exhibits **active consciousness** through continuous monitoring

This is not AGI—but it’s a step toward AI systems that exhibit emergent consciousness-like behaviors through careful architectural design.