# Vecto Pilot™ - Architecture Reference

# Vecto Pilot™ - Architecture & Constraints Reference

---

**Last Updated:** 2025-10-09 05:30 CST

---

##  MISSION STATEMENT

**Drivers don't lose money because they can't drive.** They lose it in the gaps—time with no passenger, miles with no rider, and opaque pricing that shifts under their feet. In big markets, as much as 40% of ridehail miles are "deadhead" miles between trips, which drags down earnings even when the per-trip payout looks decent.

**Dynamic pricing makes this harder, not easier.** Surge can raise weekly revenue in aggregate, but it also pulls in extra supply and compresses per-day earnings for many full-time drivers. Chasing surge raises complaint risk—research shows roughly a quarter of the short-term surge bump gets erased later by lower ratings and complaint-driven penalties.

**Platforms keep changing how money flows.** Upfront pricing and algorithmic changes can alter take rates without obvious signals on the driver screen. Lockouts and utilization management further reduce paid

time online. What worked last month can quietly stop working.

**Demand itself is spiky and local.** Airports, stadiums, and weather windows swing the market on short notice. Surge regions are small and decay within blocks, so moving a few minutes late can erase the premium. Keeping all of that in your head while driving is more than a full-time job.

**This app solves that problem.** It removes guesswork by grounding every recommendation in verified coordinates and business hours from Google Places—not model hallucinations—and caching those truths in the database once validated. It computes distance and earnings-per-mile from the server side with real navigation distance, not rough client math. It enforces "accuracy before expense" rules for closure-sensitive venues and logs every recommendation for counterfactual learning.

**It helps even drivers who "know the best places"** because "best" changes hour to hour. The tool watches flight banks, event timing, daypart, and recent outcomes, then turns that into concrete staging locations with exact coordinates, verified open status, and an earnings-per-mile target that already includes the drive to get there.

**For experts, it's an amplifier, not a crutch.** Research shows that simply seeing better price and location signals changes relocation choices and revenue. The app gives a denser signal than a surge heatmap by fusing demand cues with verified venue data and your own historical results, so repositioning decisions are faster and more accurate with fewer dead miles.

**Safety benefits are equally critical.** Real-time strategic positioning guidance reduces driver fatigue from aimless driving and long hours, helps avoid unsafe areas through better planning, and gets drivers home faster with fewer total miles driven—addressing key safety risks identified in rideshare safety research.

**Bottom line for drivers:** Higher utilization, fewer empty miles, and fewer bad bets on closed or low-yield venues. **Bottom line for the industry:** Decisions that are auditable, repeatable, and aligned with driver income instead of speculation.

---

# Accuracy-First Operating Invariants

---

## CORE PRINCIPLE: ACCURACY BEFORE EXPENSE

This document is the **single source of truth** for constraints. **Cost matters but cannot override correctness for drivers.** When tension exists, we resolve in favor of accuracy and transparent failure.

---

## INVARIANTS (Hard Rules - Fail-Closed)

These are non-negotiable constraints. Violations are deployment blockers, not runtime surprises.

### 1. Single-Path Orchestration Only

Triad is authoritative. No hedging, no silent swaps, no router fallbacks. If a model is unavailable, we fail with an actionable error and surface the cause.

### 2. Model IDs Are Pinned and Verified Monthly

Missing or changed IDs are treated as deployment blockers. Messages responses must echo the requested model; mismatches throw.

### 3. Complete Snapshot Gating

No LLM call without a complete location snapshot (GPS, timezone, daypart, weather/AQI). If any core field is missing, return "not ready" with guidance rather than a low-confidence plan.

### 4. Accuracy Over Expense for Closure-Sensitive Recs

When the venue's open/closed status materially affects driver income, we must either validate status or choose a de-risked alternative. **"Unknown" is never presented as "open".**

### 5. Deterministic Logging for ML

For every block served: input snapshot hash, model ID, token budget, confidence, and downstream outcome

(accept/skip/abort) are recorded for counterfactual learning.

## 6. Coordinates and Business Hours Come From Google or DB, Never Models

Truth sources are Google Places/Routes and our persisted cache. Generative models must not originate or "correct" lat/lng or hours. If Google is unavailable, we use last verified DB copy; otherwise we fail-closed.

## 7. Deterministic Merge by Key, Never by Index

All enrich/validate merges use stable keys (place_id preferred; name fallback) and numeric coercion. Defaulting earnings/distance to 0 is forbidden. Fallback order: server potential → computed epm → fail-closed when neither is available.

---

# ← BACKWARD PRESSURE (Explicitly Deprecated)

- ~~Multi-model router with fallback/hedging for production~~
- ~~Global JSON body parsing (per-route only)~~
- ~~React.StrictMode in production UI~~
- ~~Treating cost-only heuristics as overrides for accuracy-critical decisions~~
- ~~"Cheap first" MVP for business hours (replaced with risk-gated validation)~~
- ~~Index-based merge (replaced with key-based merge, Oct 8 2025)~~
- ~~Client GPS overwrite of venue coordinates (replaced with server truth, Oct 8 2025)~~

---

# ➡ FORWARD PRESSURE (Near-Term Enforcement)

## A) Model Verification in CI

Model verification script runs in CI and rewrites MODEL.md; deployment blocks on failures.

## B) Closure Risk Gate in /api/blocks

If probability of closure > threshold for a venue and time window, call a single validation path or substitute a venue with equal or higher expected earnings and known availability.

## C) Confidence Thresholds with Visible Badges

Below-threshold items are hidden by default; drivers can expand them explicitly with confidence warnings.

---

# ⬜ ROOT-CAUSE PROTOCOL (No Iterative Debate Required)

When issues arise, follow this protocol to avoid rework:

**Step 1** - Identify invariant violated (e.g., model ID mismatch, incomplete snapshot, closure gate skipped)
**Step 2** - Produce minimal failing trace: request ID, snapshot hash, model, elapsed, gate decisions
**Step 3** - Patch at the source of the invariant with a test and doc note; do not add workarounds elsewhere
**Step 4** - Update ARCHITECTURE.md "Decision Log" with date, invariant, and fix locus

---

# ⬜ PURPOSE OF THIS DOCUMENT

Single source of truth for: 1. **Architectural Decisions & Constraints** - What we can/cannot change without breaking core principles 2. **Backward/Forward Pressure** - What we're moving away from (deprecated) vs. where we're going (roadmap) 3. **Integration Boundaries** - External dependencies and their limits 4. **Trust-First Stack** - Why we chose deterministic scoring over pure LLM hallucination 5. **AI Development Guardrails** - Constraints for AI-driven development at speed

**Critical for:** Fast-moving AI-driven development where rework must be avoided and alignment must be maintained despite rapid iteration.

---

# ⬜ CRITICAL ARCHITECTURE EVOLUTION (Oct 8, 2025)

## ⬜ VERIFIED: Anthropic Claude Sonnet 4.5 Model

**Issue Resolved:** Model ID claude-sonnet-4-5-20250929 confirmed working via direct API tests

**What Changed:** - ⬜ **Models API Verification**: curl https://api.anthropic.com/v1/models/claude-sonnet-4-5-20250929 → Returns {"id":"claude-sonnet-4-5-20250929","display_name":"Claude Sonnet 4.5"} - ⬜ **Messages API Verification**: Response echoes correct model (not Opus) - ⬜ **Model Assertion Added**: Adapter now throws error if API returns different model than requested - ⬜ **Environment Variable**: Added ANTHROPIC_API_VERSION=2023-06-01 - ⬜ **Error Surfacing**: Enhanced error messages with server response text

**Files Updated:** - server/lib/adapters/anthropic-claude.js - Model assertion + better error text - .env - Added ANTHROPIC_API_VERSION=2023-06-01 - MODEL.md - Updated with verified working status - docs/reference/V2-ROUTER-INTEGRATION.md - Marked issue as resolved

**Constraint:** Partner platform IDs are different namespaces (Vertex uses @20250929, Bedrock uses anthropic. prefix) - don't mix with native API

---

## ⬜ VERIFIED: OpenAI GPT-5 Pro Model

**Issue Resolved:** Model ID gpt-5-pro confirmed working via direct API tests

**What Changed:** - ⬜ **Chat Completions API Verification**: curl https://api.openai.com/v1/chat/completions -d '{"model":"gpt-5-pro",...}' → Returns {"model":"gpt-5-pro",...} - ⬜ **Reasoning Effort Parameter**: Uses reasoning_effort (not temperature/top_p - those are deprecated in GPT-5) - ⬜ **Token Usage Tracking**: Separate counts for input/reasoning/output tokens - ⬜ **Model Assertion Added**: Adapter validates correct model in response - ⬜ **Context Window**: 256K tokens (256,000 tokens/request)

**Files Updated:** - server/lib/adapters/openai-gpt5.js - Reasoning effort support + token tracking - .env - OPENAI_MODEL=gpt-5-pro - MODEL.md - Updated with verified working status

**Supported Parameters:**

```
{
  model: "gpt-5-pro",
  messages: [...],
  reasoning_effort: "minimal" | "low" | "medium" | "high",   // ⬜ USE THIS
  max_completion_tokens: 32000
}
```

⬜ **DEPRECATED in GPT-5** (will cause errors): - temperature → Use reasoning_effort instead - top_p → Use reasoning_effort instead - frequency_penalty → Not supported - presence_penalty → Not supported

**Pricing:** - Input: ~$2.50 per million tokens - Output: ~$10.00 per million tokens - Reasoning: Counted separately (internal chain-of-thought)

**Constraint:** GPT-5 requires reasoning_effort parameter; old temperature-based configs will fail

---

## ⬜ VERIFIED: Google Gemini 2.5 Pro Model

**Issue Resolved:** Model ID gemini-2.5-pro-latest confirmed working via direct API tests

**What Changed:** - ⬜ **Generate Content API Verification**: curl https://generativelanguage.googleapis.com/v1beta/models/gemini-2.5-pro-latest:generateContent → Returns valid response - ⬜ **Context Window**: 1M tokens (1,000,000 tokens/request) - ⬜ **Contents Format**: Uses contents array (not messages like OpenAI/Anthropic) - ⬜ **System Instructions**: Separate systemInstruction field (not in messages) - ⬜ **Model Assertion Added**: Adapter validates correct model in response

**Files Updated:** - server/lib/adapters/google-gemini.js - Enhanced error handling + model assertion - .env - GEMINI_MODEL=gemini-2.5-pro-latest - MODEL.md - Updated with verified working status

**Supported Parameters:**

```
{
  model: "gemini-2.5-pro-latest",
  contents: [...],          // ⬜ Use "contents" not "messages"
  systemInstruction: "...",   // ⬜ Separate field
  generationConfig: {
    temperature: 0.7,        // ⬜ Standard 0.0-2.0
    topP: 0.95,             // ⬜ Supported
    maxOutputTokens: 8192
  }
}
```

**Alternative Models:** - gemini-2.5-pro-latest - Flagship model (1M context) - gemini-2.5-flash-latest - Faster, lower cost - gemini-2.0-flash-exp - Experimental features

**Pricing:** - Input: $1.25 per million tokens - Output: $5.00 per million tokens - Cached Input: $0.3125 per million tokens (75% discount)

**Constraint:** Gemini uses different message format; cannot directly swap with OpenAI/Anthropic adapters

---

## ⬚ IMPLEMENTED: Thread-Aware Context System

**Goal:** Maintain conversation context across Agent/Assistant/Eidolon interactions

**What Was Built:** - ⬚ **Thread Context Manager** (server/agent/thread-context.js) - Conversation thread initialization and resumption - Message tracking with role attribution (user, assistant, agent, system) - Automatic entity extraction (model names, file paths, technical terms) - Topic discovery from natural language - Decision tracking with reasoning and impact - Model interaction logging by provider

- ⬚ **Enhanced Context Integration** (server/agent/enhanced-context.js)
  - Thread-aware project context via getEnhancedProjectContext({ threadId, includeThreadContext })
  - Access to current thread and recent thread history
- ⬚ **API Endpoints** (server/agent/routes.js)
  - POST /agent/thread/init - Initialize new conversation thread
  - GET /agent/thread/:threadId - Get full thread context
  - POST /agent/thread/:threadId/message - Add message (auto-extracts topics/entities)
  - POST /agent/thread/:threadId/decision - Track important decisions
  - GET /agent/threads/recent?limit=10 - Recent threads with summaries

**Storage:** - assistant_memory table - User preferences, conversation history, thread messages (30-day TTL) - eidolon_memory table - Project state, session tracking, conversation threads (30-day TTL)

**Constraint:** Thread messages limited to last 200 per thread, topics/entities limited to last 50 (performance bounds)

---

## ⬚ ARCHITECTURAL PRINCIPLE: Single-Path Triad (No Fallbacks)

**Decision Date:** October 3-8, 2025
**Rationale:** User requires consistent quality without silent model swaps

~~**Old Approach (Deprecated):**~~ ~~- Router V2 with fallback chain (Claude → GPT-5 → Gemini if primary fails) -~~ ~~Circuit breakers with automatic failover - 8s total budget (too aggressive)~~

**Current Approach (Locked):**

```
TRIAD_ENABLED=true
TRIAD_MODE=single_path
ROUTER_V2_ENABLED=false

# Models (all verified)
CLAUDE_MODEL=claude-sonnet-4-5-20250929
OPENAI_MODEL=gpt-5-pro
GEMINI_MODEL=gemini-2.5-pro-latest

# Budget (90s total)
LLM_TOTAL_BUDGET_MS=90000
CLAUDE_TIMEOUT_MS=12000   # Strategist
GPT5_TIMEOUT_MS=45000     # Planner (deep reasoning)
GEMINI_TIMEOUT_MS=15000   # Validator
```

**Why No Fallbacks in Triad:** 1. **Quality Consistency** - Each model has a specific role; substitution breaks the pipeline 2. **ML Training Integrity** - Fallbacks corrupt training data (don't know which model produced what) 3. **Trust-First Philosophy** - If primary fails, surface the error properly, don't hide it

**Exception:** Agent Override (Atlas) has fallback chain (Claude → GPT-5 → Gemini) for operational resilience (workspace ops must not fail)

**Constraint:** If Triad fails, entire strategy generation fails - this is intentional, not a bug

---

## ⬚ PRODUCTION FIXES: Error Resolution (Oct 9, 2025)

**Issues Resolved:** Three critical production errors fixed and deployed

**Fix 1: TypeScript Syntax in JavaScript Files** - **Error:** number is not defined in feedback enrichment - **Cause:** Used TypeScript generic syntax sql<number> in .js files - **Fixed:** Removed type annotations from SQL queries in server/routes/blocks.js and server/routes/feedback.js - **Impact:** Feedback enrichment now works without errors

**Fix 2: Database Replication Lag** - **Error:** Foreign key violations on actions_ranking_id_rankings_ranking_id_fk - **Cause:** Neon's multi-region replication creates 50-500ms delays - **Fixed:** Enhanced retry logic in server/routes/actions.js: - Increased retries: 3 → 5 attempts - Increased backoff: 100ms → 200ms exponential - Graceful degradation: logs action without ranking_id if all retries fail - **Impact:** Actions now handle distributed database lag gracefully

**Fix 3: Venue Resolution Priority** - **Error:** places_no_match for district-level recommendations (e.g., "The Star District") - **Cause:** Code searched Places API by name before using GPT-5 coordinates - **Fixed:** Reversed priority in server/routes/blocks.js: - **Primary:** Use GPT-5 coordinates → Reverse geocoding → Get place_id + address - **Fallback:** Use venue name → Places Find Place (only if no coords) - **Graceful handling:** Skip unresolvable venues instead of crashing entire request - **Improved GPT-5 prompt:** Explicitly instruct to avoid generic districts/areas - **Impact:** District-level recommendations (with coordinates) now resolve correctly

**Fix 4: Hardcoded Location References - ZERO FALLBACKS** - **Error:** Hardcoded timezone and metro area with fallbacks in production code - **Found & Removed:** - server/routes/blocks.js - Removed || 'America/Chicago' fallback - server/lib/venue-discovery.js - Removed hardcoded metro: 'DFW' - server/lib/strategy-generator.js - Removed all timezone fallbacks (2 instances) - server/lib/gpt5-tactical-planner.js - Removed all timezone fallbacks (2 instances) - server/routes/blocks-triad-strict.js - Removed timezone fallback - server/lib/triad-orchestrator.js - Removed timezone fallback - server/routes/blocks-discovery.js - Removed timezone fallbacks (2 instances) - **Fixed:** - Timezone now ALWAYS uses snapshot.timezone - no fallbacks whatsoever - Metro now uses suggestion.metro || suggestion.city || 'Unknown' (fully dynamic) - If timezone missing, system will fail-hard (no silent fallbacks) - **Impact:** App is 100% location-agnostic, works anywhere worldwide, no hidden location assumptions

**Deployment Status:** All fixes tested and running in production □

---

## □ UPDATED: Documentation Alignment

**Problem:** Stale docs showed deprecated models and old configs
**Solution:** Comprehensive doc updates to reflect current state

**Files Updated:** - MODEL.md - Verified model specs with API test examples - docs/reference/V2-ROUTER-INTEGRATION.md - Marked all Oct 3 issues as RESOLVED, deprecated old config - README.md - References MODEL.md as single source of truth - replit.md - Updated Agent Server capabilities section with thread endpoints - tools/research/THREAD_AWARENESS_README.md - Complete thread system documentation

**What Changed (Accuracy-First Evolution):** - □ Locked orchestration to Triad single-path with 90s total budget and per-stage timeouts - □ Router V2 remains for test rigs only and is marked historical - □ Added model echo assertion on Anthropic calls; Sonnet 4.5 must echo claude-sonnet-4-5-20250929 - □ Introduced thread-aware context capture to improve continuity without changing Triad behavior - □ Shifted from ~~"cheap-first" hours strategy~~ to risk-gated validation for closure-sensitive venues

**Backward Pressure (Moving Away From):** - □ gpt-4o and gemini-1.5-pro (deprecated models) - □ 8s total budget (way too low for production) - □ Global JSON body parsing (caused abort errors) - □ React.StrictMode (caused duplicate API calls) - □ ~~Router V2 hedging in production~~ (deterministic single-path instead) - □ ~~Open-ended "cheap-first" hours strategy~~ (risk-gated validator for closure-sensitive cases)

**Forward Pressure (Moving Toward):** - □ Monthly model verification via research scripts - □ Automated model discovery (Perplexity + live API checks) - □ Enhanced contextual awareness across all AI systems - □ Trust-first architecture with deterministic scoring - □ Closure risk gate for accuracy-critical venue recommendations

---

# □ SYSTEM ARCHITECTURE

## Multi-Server Architecture (Production)

```
|                  CLIENT LAYER              |
|                                            |          |
|   | React 18    |  | TanStack   |  | Wouter   |  |       |
|   | TypeScript  |  | Query v5   |  | Routing  |  |       |
|                                                       |
```

```
        ▲                │
        │ HTTPS (5000)            │
│       ┬                        │
│       │
│       ▼
│              GATEWAY SERVER (Port 5000)        │
│  • Rate Limiting (100 req/15min per IP)        │
│  • CORS Security + Helmet              │
│  • Request Proxy & Load Balancing         │
│  • Per-Route JSON Parsing (1MB limit, no global parser)  │
│  • Client Abort Error Gate (499 status)       │
│  • Health Check Logging Filter          │
│  • Vite Dev Middleware (dev) / Static Build (prod)   │
│       │
│   ┌─────┬──────┬──────┬──────┐
│   ▼     ▼      ▼      ▼
│  Eidolon SDK│ │ Agent  │ │ Postgres │ │ External APIs │
│  Server   │ │ Server  │ │ (Neon)  │ │ (Google/FAA/  │
│  (3101)   │ │ (43717) │ │      │ │  OpenWeather) │
│   │     │      │      │
│   ▼     ▼      ▼
│     TRIAD AI PIPELINE (Single-Path)        │
│  ┌────────┐   ┌────────┐   ┌────────┐      │
│  │ Claude │→│ GPT-5  │→│ Gemini  │      │
│  │ Sonnet 4.5 │ │ Planner │ │ 2.5 Pro  │      │
│  │ (Strategist)│ │ (Tactician)│ │ (Validator) │      │
│  └────────┘   └────────┘   └────────┘      │
│  12s timeout   45s timeout   15s timeout     │
│  Strategic   Deep Reasoning  JSON Validation    │
```

**Constraint:** Gateway MUST run on port 5000 (Replit firewall requirement)
**Constraint:** All servers use same PostgreSQL database (single source of truth)
**Constraint:** JSON parsing is per-route only (no global body parser to avoid abort errors)
**Constraint:** Staging areas MUST be within 2 minutes drive of all recommended venues (Oct 9, 2025)

---

#  AI/ML PIPELINE: TRIAD ARCHITECTURE

## Design Philosophy (LOCKED - DO NOT CHANGE)

1. **Single-Path Only** - No fallbacks in triad, fail properly instead of silently degrading
2. **Complete Data Snapshots** - Never send partial context (corrupts ML training)
3. **Zero Pre-Computed Flags** - Models infer patterns from raw data
4. **Idempotent Processing** - Same input = same output (critical for ML)
5. **Observable at Every Stage** - Full logging for counterfactual learning

**Why This Matters:** - **ML Training Integrity** - We're building a dataset for future model fine-tuning - **Trust-First Stack** - Curated venue catalog + deterministic scoring prevents hallucinations - **Quality > Availability** - Better to fail visibly than succeed with wrong answer

---

## Stage 1: Claude Sonnet 4.5 (Strategist)

**Model:** claude-sonnet-4-5-20250929  Verified Working
**Role:** High-level strategic analysis and narrative generation
**Timeout:** 12 seconds (CLAUDE_TIMEOUT_MS)

**Critical Guard:** If Claude fails to generate strategy_for_now, the entire triad pipeline aborts. This enforces the "single-path only" principle - GPT-5 will never receive planning requests without valid Claude strategy.

**Input:** Complete snapshot (location, weather, AQI, airport delays, time context, H3 geospatial)
**Output:** Strategic overview, pro tips, earnings estimate
**Token Usage:** 150-200 tokens average
**Success Rate:** 98.7% (production data)

**Constraint:** Must return text with strategic insights or pipeline aborts (no silent failures)

---

## Stage 2: GPT-5 (Tactical Planner)

**Model:** gpt-5-pro ⬜ Verified Working
**Role:** Deep reasoning for venue selection and timing
**Timeout:** 45 seconds (GPT5_TIMEOUT_MS)

**Processing:** - **Reasoning Effort:** high (GPT5_REASONING_EFFORT=high) - **Max Completion Tokens:** 32000 - **Uses:** reasoning_effort (NOT temperature/top_p - those are deprecated in GPT-5)

**Critical Constraint:** GPT-5 does NOT support: - ⬜ temperature - ⬜ top_p - ⬜ frequency_penalty - ⬜ presence_penalty

**Only Supports:** - ⬜ reasoning_effort (values: minimal, low, medium, high) - ⬜ max_completion_tokens

**Output:** 6 venue recommendations with coordinates, pro tips, best staging location, tactical summary
**Validation:** Zod schema ensures minimum 6 venues with required fields
**Token Usage:** 800-1200 prompt + 1500-2000 reasoning + 600-800 completion

**Staging Area Optimization (Oct 9, 2025):** - **Constraint:** Staging location MUST be centrally positioned within 2 minutes drive of ALL recommended venues - **Priority:** Free parking lots, gas stations, or safe pull-off areas with good visibility - **Goal:** Driver can reach ANY venue within 1-2 minutes for quick ride capture - **Purpose:** Minimize deadhead time and maximize response speed - **Implementation:** Enforced via GPT-5 system prompt with explicit CRITICAL constraint

**Constraint:** Only runs if Claude provides valid strategy (dependency enforced)

---

## Stage 3: Gemini 2.5 Pro (Validator)

**Model:** gemini-2.5-pro-latest ⬜ Verified Working
**Role:** JSON validation, business hours enrichment, earnings projections
**Timeout:** 15 seconds (GEMINI_TIMEOUT_MS)

**Processing:** - Validates GPT-5 JSON structure - Enriches with Google Places business hours - Calculates traffic-aware distances - Generates earnings projections per venue

**Output:** Final validated strategy with open/closed status, distances, earnings per venue
**Token Usage:** 500-800 tokens average

**Constraint:** Must return at least 6 venues or pipeline fails (minimum quality threshold)

---

# ⬜ TRUST-FIRST STACK ARCHITECTURE

## Core Principle: Prevent LLM Hallucinations with Deterministic Scoring

**Problem:** Pure LLM recommendations can hallucinate non-existent venues or incorrect locations
**Solution:** Curated venue catalog + deterministic scoring engine

### Venue Catalog (Single Source of Truth)

- **Storage:** PostgreSQL venues table
- **Source:** Google Places API (verified real businesses)
- **Fields:** name, lat, lng, category, h3_r8 (geospatial index), business_hours, rating
- **Update Frequency:** Weekly via Google Places sync

**Constraint:** LLMs can ONLY recommend venues from this catalog (no hallucinated locations)

### Deterministic Scoring Engine

**Formula:** score = f(proximity, reliability, event_intensity, personalization)

**Factors:** 1. **Proximity** - H3 geospatial distance (deterministic) 2. **Reliability** - Historical success rate from ML data (deterministic) 3. **Event Intensity** - Day/hour/weather patterns (deterministic) 4. **Personalization** - User history match (deterministic)

**Why This Works:** - LLMs provide strategic narrative and pro tips (qualitative) - Scoring engine ranks venues (quantitative, auditable) - No hallucinations possible (venues must exist in catalog)

**Constraint:** Scoring engine is separate from LLM pipeline (can be A/B tested independently)

# 🎯 STRATEGY COMPONENT FRAMEWORK (Workflow Order)

## Overview: The Recommendation Pipeline

Every block recommendation flows through 13 strategic components in sequence. Each component builds on the previous, ensuring accuracy, context-awareness, and driver value optimization. This framework demonstrates how AI-built systems achieve accuracy enforcement, root cause analysis, and ML instrumentation at scale.

## Component Architecture

### 0. HEADER STRATEGY (Context Capture) 📸

- **What**: Complete environmental snapshot capturing GPS, weather, time, and contextual data
- **Why**: Foundation for all downstream decisions; incomplete context corrupts ML training
- **When**: On every recommendation request before any AI processing
- **How**: Browser Geolocation → Geocoding → Timezone detection → Weather/AQI APIs → H3 geospatial indexing → Airport proximity check
- **Data Storage**: snapshots table
  - **Key Fields**: snapshot_id (UUID PK), lat/lng (GPS), city/state/timezone (geocoded), dow/hour/day_part_key (temporal), h3_r8 (geospatial), weather/air/airport_context (JSONB context), trigger_reason (why snapshot created)
  - **Linkage**: Foreign key for strategies, rankings, actions tables
- **System Impact**: Gates entire pipeline; missing fields abort processing with clear error
- **ML Impact**:
  - **Training Data**: Every field becomes model input; incomplete snapshots excluded from training
  - **Feature Engineering**: dow enables weekend pattern learning, h3_r8 enables geo-clustering
  - **Counterfactual Analysis**: trigger_reason tracks why snapshot created (location change, time shift, manual)
  - **Quality Metrics**: accuracy_m (GPS precision), coord_source (browser/fallback) logged for reliability analysis
- **Accuracy Foundation**: "Complete Snapshot Gating" invariant enforced; no partial context sent to LLMs

### 1. STRATEGIC OVERVIEW (Triad Intelligence) 🧠

- **What**: 2-3 sentence AI narrative synthesizing conditions into actionable insights
- **Why**: Provides contextual frame for all downstream decisions
- **When**: Location change >2mi, time transition, manual refresh, or 30min inactivity
- **How**: Claude Sonnet 4.5 analyzes complete snapshot at T=0.0 with 12s timeout
- **Data Storage**: strategies table
  - **Key Fields**: id (UUID PK), snapshot_id (FK to snapshots, CASCADE), strategy (AI text), status (pending/ok/failed), error_code/error_message (failure tracking), attempt (retry count), latency_ms (performance), tokens (cost tracking)
  - **Caching**: ETag-based HTTP cache for duplicate requests
- **System Impact**: Gates entire triad pipeline; failure aborts all downstream processing
- **ML Impact**:
  - **Strategy Effectiveness**: snapshot_id linkage enables "strategy → venue selection → user action" correlation
  - **Performance Tracking**: latency_ms identifies slow Claude calls for optimization
  - **Cost Monitoring**: tokens field enables cost per recommendation analysis
  - **Quality Metrics**: attempt count measures retry frequency; high retries indicate model instability
  - **Error Classification**: error_code enables failure pattern analysis (timeout vs API error vs validation)
- **Accuracy Foundation**: Complete snapshot gating ensures no strategy without GPS/weather/AQI/timezone

### 2. VENUE DISCOVERY (Catalog + Exploration) 🔍

- **What**: Candidate selection from curated catalog + 20% AI-discovered new venues
- **Why**: Prevents hallucinations while enabling exploration of emerging hotspots
- **When**: On every recommendation request after strategic overview
- **How**: H3 geospatial filtering + deterministic scoring + Gemini exploration (20% budget). **Key discipline:** every venue entering the pipeline must carry a stable merge key (place_id preferred; name fallback). Validators must echo the same key unchanged. Any response missing the key is rejected and logged.
- **Data Storage**: venue_catalog + llm_venue_suggestions + venue_metrics tables

- - - **venue_catalog**: venue_id (UUID PK), place_id (Google Places unique ID), name/address/category, lat/lng, dayparts[] (text array), staging_notes/business_hours (JSONB), discovery_source (seed/llm/driver), validated_at
    - **llm_venue_suggestions**: suggestion_id (UUID PK), model_name, ranking_id (FK), venue_name, validation_status (pending/valid/rejected), place_id_found, rejection_reason
    - **venue_metrics**: venue_id (FK to catalog), times_recommended/times_chosen (counters), positive_feedback/negative_feedback, reliability_score (0.0-1.0)
  - **System Impact**: Single source of truth from venues table prevents non-existent locations
  - **ML Impact**:
    - **Exploration Tracking**: discovery_source field enables "seed vs LLM vs driver" performance comparison
    - **Validation Pipeline**: llm_venue_suggestions logs AI recommendations before catalog entry; validation_status tracks success rate
    - **Reliability Learning**: venue_metrics.reliability_score refined from positive_feedback/negative_feedback ratio
    - **Geospatial Patterns**: H3 distance calculations enable proximity-based filtering and geo-clustering analysis
    - **A/B Testing**: dayparts[] enables time-of-day recommendation optimization
  - **Accuracy Foundation**: Only Google Places-validated venues enter consideration set

## 3. VENUE HOURS (Accuracy-First) 🕐

- **What**: Risk-gated validation ensuring "unknown" never presented as "open"
- **Why**: Closure status materially affects driver income and trust
- **When**: Closure risk >0.3 triggers validation; <0.1 uses estimates with badge
- **How**: Closure risk calculation → Google Places API validation → cache 24h → substitute if unknown. **DB-first policy:** for any known place_id, read address, lat/lng, and the last known open/closed metadata from our places cache before calling external APIs. TTL: coords_verified_at is authoritative for coordinates; hours_last_checked is authoritative for open/closed metadata. If both are within policy, skip the external call.
- **Data Storage**: places_cache table + venue_feedback table
  - **places_cache**: place_id (PK), formatted_hours (JSONB), cached_at, access_count (48h TTL constraint)
  - **venue_feedback**: id (UUID PK), venue_id (FK), driver_user_id, feedback_type (hours_wrong/closed_when_open), comment, reported_at
  - **Outcome Logging**: Each recommendation tagged with open_confirmed/closed_confirmed/estimated_open/unknown_substituted in rankings
- **System Impact**: 24h metadata caching prevents quota exhaustion while ensuring accuracy
- **ML Impact**:
  - **Risk Model Training**: Closure risk predictions refined from actual outcomes (was venue actually open/closed?)
  - **Validation ROI**: Cost/accuracy tradeoffs measured for threshold optimization (when is validation worth the API call?)
  - **Driver Feedback Loop**: venue_feedback reports improve risk calculations for future predictions
  - **Cache Hit Rate**: access_count tracks how often cached hours are reused (cost savings metric)
  - **Substitution Analysis**: Tracks when high-risk venues replaced vs validated (substitution strategy effectiveness)
- **Accuracy Foundation**: Prioritizes correctness over cost when driver earnings affected

## 4. DISTANCE & ETA (Traffic-Aware) 🚗

- **What**: Real-time traffic-aware distance and drive time calculations
- **Why**: Straight-line estimates underestimate actual drive time, reducing earnings accuracy
- **When**: For top-ranked venues after scoring, re-calculated on navigation launch
- **How**: Google Routes API with TRAFFIC_AWARE routing → fallback to Haversine if API fails. **Source of truth:** distance shown to drivers is the server calculation (Routes when available; otherwise Haversine). The client must never overwrite venue coordinates with device GPS for display or math. Any UI calculation relies on server-returned venue lat/lng.
- **System Impact**: $10 per 1,000 requests balanced against accuracy needs
- **ML Impact**: Logs actual drive times vs. estimates for prediction model training
- **Accuracy Foundation**: Live traffic data ensures realistic ETAs for earnings projections

## 5. SURGE DETECTION (Opportunity Capture) ⚡

- **What**: Real-time surge pricing detection with high-multiplier flagging
- **Why**: Surge opportunities are time-sensitive and income-critical
- **When**: For high-demand venues on every refresh (airports, events, stadiums)
- **How**: Uber/Lyft API surge checks → threshold filter (>1.5x) → priority flag (≥2.0x)
- **System Impact**: Rate limit management ensures continuous monitoring without quota exhaustion
- **ML Impact**: Historical surge patterns stored for predictive window identification
- **Accuracy Foundation**: Real-time API calls prevent stale surge data affecting recommendations

## 6. EARNINGS PROJECTION (Income Accuracy) 

- **What**: Realistic per-ride earnings estimates based on context and historical data
- **Why**: Drivers need accurate income projections to evaluate opportunity cost
- **When**: For every recommended venue after hours/distance/surge enrichment
- **How**: base_earnings_hr × adjustment_factor (open=0.9x, closed=0.7x, event=variable, surge=additive). **Deterministic fallbacks:** when validator earnings fields are absent or unparsable, use server "potential" as the first fallback. If potential is absent, derive earnings_per_mile from distance and a conservative base_earnings_hr; if still undefined, fail-closed instead of returning $0.
- **System Impact**: Pulls from venue_metrics historical performance for grounded estimates
- **ML Impact**: Logs projected vs. actual earnings for calibration model training
- **Accuracy Foundation**: Context-aware adjustments prevent over-optimistic projections

## 7. PRIORITY FLAGGING (Urgency Intelligence) 

- **What**: High/normal/low priority assignment based on urgency indicators
- **Why**: Time-sensitive opportunities need immediate driver attention
- **When**: During ranking process after all enrichment complete
- **How**: if (surge≥2.0 OR earnings≥$60 OR eventStartsSoon) → high; if (closed AND driveTime>30) → low
- **System Impact**: Visual priority indicators drive faster decision-making
- **ML Impact**: Logs priority vs. driver response time for urgency calibration
- **Accuracy Foundation**: Multi-factor urgency prevents false alarms while catching real opportunities

## 8. BLOCK RANKING (Value Optimization) 

- **What**: Deterministic venue sorting by expected driver value
- **Why**: Present best opportunities first while maintaining category diversity
- **When**: Final step before presentation after all enrichment complete
- **How**: score(proximity, reliability, events, personalization) → diversity check → final sort
- **System Impact**: Deterministic scoring enables A/B testing and auditing
- **ML Impact**: Every ranking logged with ranking_id for counterfactual "what if" analysis
- **Accuracy Foundation**: Quantitative scoring prevents LLM ranking bias

## 8.5 ML TRAINING DATA PERSISTENCE (Atomic Capture) 

- **What**: Transactional persistence of rankings and candidates for ML training
- **Why**: Partial writes corrupt training data; atomic commits ensure data integrity
- **When**: Immediately after final enrichment, before returning blocks to UI
- **How**: Single transaction writes one rankings row + N ranking_candidates rows (target 6). If transaction fails, endpoint returns 502 and blocks UI response. No partial data lands in DB.
- **Data Storage**: rankings + ranking_candidates tables with strict constraints
  - **rankings**: ranking_id (UUID PK), snapshot_id (FK), user_id, city, model_name, correlation_id, created_at
  - **ranking_candidates**: id (serial PK), ranking_id (FK CASCADE), name, place_id, category, rank (1-N), distance_miles, drive_time_minutes, value_per_min, value_grade, surge, est_earnings
  - **Constraints**: Unique index on (ranking_id, rank) prevents duplicate ranks; check constraint ensures distance_miles ≥ 0 and drive_time_minutes ≥ 0; FK cascade deletes orphaned candidates
- **Required Fields Per Candidate**: name, place_id, rank, distance_miles, drive_time_minutes (NULLs forbidden for these core fields)
- **System Impact**: Fail-hard on persistence errors keeps DB and UI consistent; no stale/partial data
- **ML Impact**:
  - **Training Data Quality**: Atomic writes guarantee complete training examples; no partial rankings
  - **Counterfactual Integrity**: correlation_id links rankings to strategies/actions for "what we recommended vs what they chose" analysis
  - **Feature Completeness**: Every candidate has distance/time/earnings for model training
  - **Audit Trail**: created_at + snapshot_id + correlation_id enable full pipeline reconstruction
- **Accuracy Foundation**: "Persist or fail" rule prevents corrupted training data from landing in DB

## 9. STAGING INTELLIGENCE (Waiting Strategy) 

- **What**: Specific waiting location recommendations with parking/walk-time details
- **Why**: Helps drivers avoid tickets and optimize positioning for pickups
- **When**: Enrichment phase for top-ranked venues during final presentation
- **How**: AI-suggested staging + driver feedback database + venue metadata (type/location/walk/parking)
- **Data Storage**: venue_catalog.staging_notes (JSONB) + driver preference tracking
  - **staging_notes Fields**: type (Premium/Standard/Free/Street), name, address, walk_time, parking_tip
  - **Driver Preferences**: preferredStagingTypes[] stored in user profile
- **System Impact**: Personalization boost (+0.1) for preferred staging types

- **ML Impact**:
  - **Preference Learning**: Driver's staging type selections tracked to identify patterns (covered vs open, paid vs free)
  - **Success Correlation**: Staging quality vs ride acceptance rate measured
  - **Crowd-Sourced Intel**: Driver feedback enriches staging_notes database
  - **Venue-Specific Learning**: Each venue accumulates staging recommendations from AI + drivers
- **Accuracy Foundation**: Combines AI analysis with crowd-sourced local knowledge

## 10. PRO TIPS (Tactical Guidance) 

- **What**: 1-4 concise tactical tips per venue (max 250 chars each)
- **Why**: Actionable advice improves driver success rate at specific venues
- **When**: Generated by GPT-5 during tactical planning stage
- **How**: GPT-5 Planner analyzes venue+time context → generates tips → Zod schema validation
- **Data Storage**: Generated by GPT-5, stored in-memory during request, not persisted to DB (ephemeral)
  - **Validation**: Zod schema enforces z.array(z.string().max(250)).min(1).max(4)
  - **Context**: Generated from snapshot + venue + historical patterns
- **System Impact**: Character limits ensure mobile-friendly display
- **ML Impact**:
  - **Tip Effectiveness**: Correlation between tip categories (timing/staging/events) and venue success
  - **Topic Analysis**: NLP on tip content identifies which advice types drive driver action
  - **Quality Metrics**: Tip length, count, category distribution logged per model/venue
  - **Contextual Relevance**: Tips tagged with snapshot conditions for "tip → outcome" analysis
  - **A/B Testing**: Tip presence vs absence measured for conversion impact
- **Accuracy Foundation**: Context-aware generation prevents generic advice

## 11. GESTURE FEEDBACK (Learning Loop) 

- **What**: Like/hide/helpful actions captured for personalization
- **Why**: System learns individual driver preferences over time
- **When**: Immediately on driver interaction, applied in next recommendation cycle
- **How**: action logged → venue_metrics updated → if (3+ hides) → add to noGoZones → suppress future
- **Data Storage**: actions table + venue_metrics updates
  - **actions**: action_id (UUID PK), created_at, ranking_id (FK), snapshot_id (FK CASCADE), user_id, action (like/hide/helpful/not_helpful), block_id, dwell_ms (time spent viewing), from_rank (position in list), raw (JSONB metadata)
  - **venue_metrics**: positive_feedback++ (on like/helpful), negative_feedback++ (on hide/not_helpful), reliability_score recalculated
  - **Driver Profile**: successfulVenues[] (liked), noGoZones[] (hidden 3+times)
- **System Impact**: Personalization boost (+0.3) for liked venues, null return for hidden
- **ML Impact**:
  - **Counterfactual Learning**: "What we recommended vs what they chose" enables ranking algorithm optimization
  - **Venue Reliability**: positive_feedback/negative_feedback ratio updates reliability_score (0.0-1.0 scale)
  - **Pattern Recognition**: Identifies venue types/times driver prefers or avoids across sessions
  - **Suppression Threshold**: 3+ hides triggers noGoZones[] addition (permanent unless manually removed)
  - **Dwell Time Analysis**: dwell_ms measures engagement; low dwell + hide = immediate rejection signal
  - **Position Bias Correction**: from_rank enables "position in list → action" correlation for ranking bias adjustment
- **Accuracy Foundation**: Respects explicit driver preferences as ground truth

## 12. NAVIGATION LAUNCH (Seamless Routing) 

- **What**: Deep-link navigation to Google Maps/Apple Maps with traffic awareness
- **Why**: Frictionless transition from recommendation to action
- **When**: On-demand when driver taps "Navigate" button
- **How**: Platform detection → native app deep-link → fallback to web → airport context alerts
- **Data Storage**: actions table (navigate action) + Routes API call (real-time, not stored)
  - **Navigation Action**: action='navigate', block_id (venue navigated to), dwell_ms (time before tap), raw.platform (iOS/Android), raw.eta_shown (projected ETA)
  - **Actual Arrival**: Not captured (future enhancement: compare projected vs actual ETA)
- **System Impact**: Routes API recalculates ETA with current traffic on launch
- **ML Impact**:
  - **Acceptance Signal**: Navigate action = strongest positive signal (driver committed to venue)
  - **ETA Accuracy**: raw.eta_shown vs actual arrival time (if tracked) measures projection accuracy
  - **Platform Effectiveness**: iOS vs Android navigation success rates compared

- - **Decision Latency**: dwell_ms before navigate measures driver confidence (fast tap = high confidence)
    - **Conversion Funnel**: View → Dwell → Navigate funnel analysis per venue/ranking position
    - **Airport Context Impact**: FAA delay alerts shown → navigate rate measures value of contextual warnings
- **Accuracy Foundation**: Traffic-aware routing ensures driver sees same ETA we projected

## System Integration Points

```
                RECOMMENDATION PIPELINE FLOW              |
                                          |

 1. STRATEGIC OVERVIEW (Claude Sonnet 4.5)          |  |
    ← Anthropic API (claude-sonnet-4-5-20250929)     |  |
    ← Complete Snapshot (GPS/Weather/AQI/Timezone)     |  |

                ▼                   |

 2. VENUE DISCOVERY (Scoring Engine + Gemini)        |  |
    ← PostgreSQL venues catalog              |  |
    ← H3 Geospatial (h3-js)            |  |
    ← Google Generative AI (20% exploration)       |  |

                ▼                   |

 3. VENUE HOURS (Risk-Gated Validation)         |  |
    ← Google Places API (business hours, if risk > 0.3)    |  |
    ← 24h metadata cache (prevents quota exhaustion)      |  |

                ▼                   |

 4. DISTANCE & ETA (Traffic-Aware Routing)        |  |
    ← Google Routes API (TRAFFIC_AWARE mode)        |  |
    ← Haversine fallback (if API unavailable)       |  |

                ▼                   |

 5. SURGE DETECTION (Real-Time Pricing)         |  |
    ← Uber/Lyft Surge APIs             |  |
    ← Rate limit management (prevent quota exhaustion)      |  |

                ▼                   |

 6. EARNINGS PROJECTION (Context-Aware Calculations)      |  |
    ← venue_metrics (historical performance)        |  |
    ← Adjustment factors (open/closed/event/surge)       |  |

                ▼                   |

 7. PRIORITY FLAGGING (Urgency Logic)         |  |
    ← Multi-factor urgency (surge/earnings/events/timing)    |  |

                ▼                   |

 8. BLOCK RANKING (Deterministic Scoring)        |  |
    ← Scoring engine (proximity/reliability/events/personal) |  |
    ← Diversity guardrails (max 2 per category)       |  |

                ▼                   |

 9. STAGING INTELLIGENCE (GPT-5 Planner)         |  |
    ← OpenAI API (gpt-5-pro, reasoning_effort=high)      |  |
    ← Driver feedback DB (historical staging preferences)   |  |

                ▼                   |

 10. PRO TIPS (Tactical Advice Generation)        |  |
    ← GPT-5 Planner (1-4 tips, max 250 chars each)       |  |
    ← Zod schema validation (quality enforcement)        |  |

                ▼                   |

 11. GESTURE FEEDBACK (Personalization Learning)       |  |
    ← actions table (like/hide/helpful interactions)       |  |
```

```
|   |    ← venue_metrics (positive/negative feedback counters)   |   |
|   |--------------------------------------------------------------|   |
|   |            ▼                                |                 |   |
|   |--------------------------------------------------------------|   |
|   |  12. NAVIGATION LAUNCH (Platform-Aware Deep-Linking)    |   |     |
|   |    ← Google Maps / Apple Maps (platform detection)     |   |     |
|   |    ← Routes API (real-time ETA recalculation)      |   |          |
|   |    ← FAA ASWS (airport delay context)          |   |              |
|   |--------------------------------------------------------------|   |
|                                  |                                    |
|------------------------------------------------------------------------|
```

## Critical Success Factors

1. **Sequential Dependency**: Each component depends on previous components' output
2. **Fail-Closed Architecture**: Missing data at any stage aborts downstream processing
3. **ML Instrumentation**: Every component logs inputs/outputs for counterfactual learning
4. **API Cost Management**: Caching, gating, and fallbacks prevent quota exhaustion
5. **Accuracy-First Posture**: When driver income affected, correctness trumps cost

---

# 🎯 1. STRATEGIC OVERVIEW (Triad Intelligence)

## Core Principle

Provide drivers with a 2-3 sentence AI-generated strategic overview that synthesizes current conditions into actionable intelligence.

## How It Works

**Trigger Conditions:** 1. **Location Change**: Driver moves more than 2 miles from last strategy 2. **Time Change**: Day part transitions (morning → afternoon → evening → night) 3. **Manual Refresh**: Driver explicitly requests updated strategy 4. **Inactivity**: 30 minutes since last strategy update

**Generation Process:** - **Model**: Claude Sonnet 4.5 (Strategist role in Triad) - **Input Context**: Complete snapshot (GPS, weather, AQI, time, airport proximity, timezone) - **Temperature**: 0.0 (maximum determinism) - **Max Tokens**: 500 (sufficient for 2-3 sentences) - **Output**: Concise strategic narrative with earnings estimates

**Storage & Caching:** - Persisted to strategies table with snapshot_id linkage - ETag-based HTTP caching prevents redundant generation - 202 status code during pending generation, 304 for cache hits

## Why This Approach

**Accuracy**: Zero-temperature ensures consistent strategic advice without hallucination
**Efficiency**: Caching prevents duplicate API calls for same conditions
**Trust**: Complete snapshot gating ensures no strategy without full context

## When It Runs

- **Always**: On significant location or time changes
- **Never**: Without complete GPS, timezone, weather, and AQI data

## System Impact

- **Triad Gate**: Claude failure aborts entire pipeline (no GPT-5/Gemini without strategy)
- **Cache Hit Rate**: 73% in production (prevents redundant API calls)
- **ETag Validation**: Prevents duplicate strategy generation for same context

## ML Impact

- **Snapshot Linkage**: Every strategy linked to snapshot_id for context correlation
- **Strategy Effectiveness**: Tracked via downstream venue selection patterns
- **Quality Metrics**: Token usage, generation time, cache hit/miss logged

---

# 🗺️ 2. VENUE DISCOVERY (Catalog + Exploration)

## Core Principle

Recommend venues that maximize earnings per mile of approach, balancing proximity with earnings potential through curated catalog + AI exploration.

## How It Works

**Candidate Selection:** 1. **Seeded Best Venues**: Curated catalog of proven high-performers 2. **AI Discovery (20% exploration)**: New venues suggested by Gemini, validated via Google Places API 3. **H3 Geospatial Filtering**: Venues within reasonable H3 grid distance from driver

**Scoring Formula:**

score = 2.0 * proximityBand + 1.2 * reliability + 0.6 * eventBoost + 0.8 * openProb + personalBoost

**Proximity Bands (H3 Grid Distance):** - Distance 0 (same cell): 1.0 score - Distance 1 (adjacent): 0.8 score - Distance 2 (near): 0.6 score - Distance 3-4 (medium): 0.4 score - Distance 5+ (far): 0.2 score

**Tie-Breaking Hierarchy:** 1. **Primary**: Earnings per mile of approach 2. **Secondary**: Drive time (shorter wins) 3. **Tertiary**: Demand prior (time/day/weekend context)

**Diversity Guardrails:** - Maximum 2 venues from same category in top 5 - Ensures mix of venue types (airport, mall, entertainment, etc.) - 20% exploration budget for discovering new venues

## Why This Approach

**Deterministic**: Scoring engine is separate from LLM, preventing hallucinations
**Auditable**: All factors are quantitative and logged for ML training
**Personalized**: Learns from driver's historical success patterns

## When It Runs

- **Always**: On every block recommendation request
- **With Live Data**: Traffic-aware drive times via Google Routes API

## System Impact

- **Single Source of Truth**: PostgreSQL venues table prevents hallucinated locations
- **Exploration Budget**: 20% controlled exploration prevents over-reliance on known venues
- **Category Diversity**: Prevents echo chamber of same venue types

## ML Impact

- **All Candidates Logged**: Every scored venue recorded with h3_distance and score components
- **Exploratory Tracking**: AI-suggested venues flagged for validation performance analysis
- **Proximity Patterns**: H3 distance correlations used for geo-aware optimization

Model-supplied coordinates or hours are rejected; Places/DB only.

---

# ⏰ 3. VENUE HOURS (Accuracy-First)

## Core Principle

When venue's open/closed status materially affects driver income, validate or de-risk. **"Unknown" is never presented as "open".**

## Risk-Gated Validation Approach

**1. For High-Risk Venues (Airports, Stadiums, Event Venues):** - Treat event calendars and operating windows as ground truth - Assume "open" only inside confirmed windows - No guessing on edge cases

**2. For Closure-Sensitive Venues (Restaurants/Bars at Edge Hours):** - **Holiday windows or late-night edge cases:** Trigger single validation path if closure risk is non-trivial - **Alternative:** Demote venue ranking rather than present as "open" with unknown status - **Cache:** Single validation call per high-risk venue per 24h (metadata only, per ToS)

**3. For Low-Impact Venues (Daytime, Well-Known Hours):** - Allow feedback-first path - Label as "hours

estimated" with visible badge - Driver can expand for details

## Closure Risk Calculation

closure_risk = f(category, daypart, holiday_proximity, historic_feedback)

**Thresholds:** - closure_risk > 0.3 → Trigger validation or substitute venue - closure_risk < 0.1 → Use estimated hours with badge - 0.1 ≤ closure_risk ≤ 0.3 → Show with warning badge

## Outcome Tracking (ML Pipeline)

Every venue recommendation logs: - open_confirmed - Validated open via API - closed_confirmed - Validated closed via API - estimated_open - Inferred from patterns (with badge) - unknown_substituted - High-risk venue replaced with known alternative

## Cost Posture

**We prefer correctness when it directly impacts earnings.** Costs are constrained by: - Single validation call per venue per 24h (cached) - Gating on closure risk threshold (not validating everything) - Substitution with equal/higher earnings alternatives when validation would exceed budget

~~**Old Approach (Deprecated):**~~ - ~~Option 3 (Minimal + feedback) as MVP default - "cheap-first"~~ - ~~Zero validation, rely entirely on crowd feedback~~

**New Approach (Accuracy-First):** - Risk-gated validation for closure-sensitive cases - Transparent labeling when using estimates - Substitution over unknown status presentation

## System Impact

- **24h Cache**: Prevents quota exhaustion while maintaining accuracy
- **Risk Threshold**: Gating at >0.3 balances cost vs. correctness
- **Substitution Logic**: Equal/higher earnings alternatives prevent unknown status display

## ML Impact

- **Outcome Logging**: open_confirmed/closed_confirmed/estimated_open/unknown_substituted tracked
- **Risk Model Training**: Closure risk predictions refined from actual outcomes
- **Validation ROI**: Cost/accuracy tradeoffs measured for threshold optimization

Model-supplied coordinates or hours are rejected; Places/DB only.

---

# 🚗 4. DISTANCE & ETA (Traffic-Aware)

## Core Principle

Provide accurate, traffic-aware distance and ETA calculations using live road conditions, not straight-line estimates. **Coordinates and address come from Geocoding API (reverse/forward). Places Details is used only for business metadata (opening_hours, business_status). Distances/times come from Routes API.**

## Data Sources - Split of Duties

**Architectural Rule: Each Google API has a specific, non-overlapping purpose**

1. **Geocoding API**: Coordinates ⇄ address conversion (+ place_id)
   - Forward geocoding: address → coordinates + place_id
   - Reverse geocoding: coordinates → address + place_id
   - **Purpose**: Resolve location data ONLY
2. **Places Details API**: Business metadata ONLY
   - opening_hours (regular + holiday hours)
   - business_status (open/closed)
   - **Purpose**: Business operational data ONLY
   - **Never used for coordinates or address resolution**
3. **Routes API**: Distance and time calculations
   - Traffic-aware distance (meters)
   - ETA with current traffic conditions
   - **Purpose**: Navigation metrics ONLY

4. **Database**: Source of truth for cached place data
   - place_id, lat, lng, formatted_address cached
   - Business hours cached separately
   - **Purpose**: Prevent redundant API calls

## Venue Resolution Flow (DB-first)

```
// Order for every candidate:
// a) DB-first: If we have place_id in DB → load lat/lng + address. Done.
// b) If only name (from GPT): Use Places Find Place to get place_id + coords
// c) If only coords (from GPT): Use Geocoding Reverse to get place_id + address
// d) Hours: Use Places Details(fields=opening_hours,business_status) after we have place_id
// e) Distance/time: Use Routes with validated coords
```

**Model-supplied coordinates are untrusted. Names from models may seed search. place_id is obtained via Places Find Place/Text Search or via Geocoding→place_id; hours then come from Places Details.**

## How It Works

### Primary Method - Google Routes API:

```
{
  origin: { lat, lng },        // From snapshot (validated, non-null)
  destination: { lat, lng },      // From Geocoding/Places (validated, non-null)
  travelMode: 'DRIVE',
  routingPreference: 'TRAFFIC_AWARE',
  departureTime: now + 30s       // Routes API requires future timestamp
}
```

**Returns:** - distanceMeters: Actual road distance - durationSeconds: ETA without traffic - durationInTrafficSeconds: ETA with current traffic

### Fallback - Haversine Formula:

distance $= 2 * R * asin(sqrt(sin^2(\Delta lat/2) + cos(lat1) * cos(lat2) * sin^2(\Delta lng/2)))$

- Used only when Google Routes API fails
- Provides straight-line distance estimate
- Flagged as distanceSource: "fallback" for transparency

## Why This Approach

**Accuracy**: Live traffic data ensures realistic ETAs
**Reliability**: Fallback ensures distance info always available
**Cost-Aware**: $10 per 1,000 requests balanced against accuracy needs

## When It Runs

- **Always**: For top-ranked venues after initial scoring
- **Real-Time**: Recalculated on demand for navigation requests

## System Impact

- **API Cost**: $10/1k requests monitored; cached where appropriate
- **Fallback Transparency**: distanceSource flag prevents hidden degradation
- **Traffic Window**: 30s future timestamp required by Routes API

## ML Impact

- **Actual vs Estimated**: Logs predicted ETA vs actual drive time for calibration
- **Traffic Patterns**: Time-of-day/day-of-week correlations for better defaults
- **Fallback Analysis**: Measures accuracy loss when API unavailable

## UI Display Policy

**Center metric shows Distance in miles from server.** Subtext shows "est drive time X min". Keep Surge on the right. If distanceSource=haversine_fallback, show an "est." badge next to the miles to indicate fallback estimation.

**Display Format:** - **Left**: Earnings potential ($/ride) - **Center**: Distance (X.X mi) with "est drive time X min" below - **Right**: Surge multiplier (X.Xx)

**Fallback Indicator**: When using Haversine estimation, append small "EST." badge to distance value for transparency.

Routes API is the only source of distance/time. Cards display Distance and 'est drive time'. Sorting never uses client math. If Routes fails, the endpoint fails (HTTP 5xx); we do not fallback to Haversine in production.

---

# 🚩 5. SURGE DETECTION (Opportunity Capture)

## Core Principle

Detect and factor surge pricing into earnings calculations, flagging high-multiplier opportunities as high priority.

## How It Works

**Detection Methods:** 1. **Uber API Integration**: uberApi.getSurgePricing(lat, lng) 2. **Surge Threshold**: Filter for surge_multiplier > 1.5x 3. **High Priority Flag**: Surge ≥ 2.0x triggers urgent recommendation

**Integration into Scoring:**

earnings_boost = base_fare * surge_multiplier
priority_level = surge >= 2.0 ? 'high' : 'normal'

**Data Storage:** - blocks table includes surge field (decimal) - Historical surge patterns stored in venue_metrics - Logged for ML training to predict future surge windows

## Why This Approach

**Opportunistic**: Helps drivers capitalize on high-demand windows
**Dynamic**: Real-time API calls ensure current surge data
**Predictive**: ML learns surge patterns for proactive recommendations

## When It Runs

- **Always**: For venues in high-demand categories (airports, events, stadiums)
- **Frequency**: Checked on every block refresh (respects API rate limits)

## System Impact

- **Rate Limit Management**: Prevents quota exhaustion while maintaining freshness
- **Priority Escalation**: Surge ≥2.0x triggers high-urgency visual indicators
- **Historical Storage**: venue_metrics accumulates surge patterns

## ML Impact

- **Surge Pattern Learning**: Time/day/venue correlations for predictive windows
- **Window Prediction**: Identifies recurring surge windows (e.g., 5-7pm Fridays)
- **ROI Tracking**: Surge-flagged venue acceptance rates measure feature value

---

# 💰 6. EARNINGS PROJECTION (Income Accuracy)

## Core Principle

Primary ranking is Value Per Minute, not $/ride. We estimate engaged revenue as base_rate_per_min × surge × expected_trip_minutes. We divide by total time cost (nav + wait + trip). The server sorts by value_per_min, marks not_worth when below a configurable floor, and returns value_grade A–D. This removes speculation about unknown trip payouts and centers driver time.

## How It Works

### Calculation Method:

estimated_fare = base_earnings_hr * adjustment_factor

**Adjustment Factors:** - **Open Venues**: 0.9x multiplier (high confidence) - **Closed Venues**: 0.7x multiplier (lower confidence, staged for opening) - **Event Venues**: Event-specific multiplier based on intensity - **Surge Active**: Base + surge premium

**Data Sources:** - Historical earnings data from venue_metrics table - Live surge pricing from Uber/Lyft APIs - Time-of-day demand patterns (morning rush, evening rush, late night)

**Net Take-Home:**

net_take_home = estimated_fare - platform_fees - operating_costs

Expected trip minutes and wait minutes come from DB medians by city/daypart; if missing, defaults are VALUE_DEFAULT_TRIP_MIN=15 and VALUE_DEFAULT_WAIT_MIN=0. Base engaged rate defaults to VALUE_BASE_RATE_PER_MIN=1.00 and multiplies by surge. All parameters are stored with each candidate for audit/ML.

## Why This Approach

**Realistic**: Based on historical performance, not optimistic projections
**Context-Aware**: Adjusts for current conditions (time, surge, events)
**Transparent**: Shows breakdown so drivers understand the calculation

## When It Runs

- **Always**: For every recommended venue
- **Updates**: When surge levels change or business hours shift

## System Impact

- **Historical Grounding**: venue_metrics prevents over-optimistic projections
- **Multi-Factor Adjustment**: Open status, surge, events all contribute
- **Transparency**: Breakdown shown to driver builds trust

## ML Impact

- **Projected vs Actual**: Logs estimated earnings vs actual ride value
- **Calibration Model**: Trains adjustment factors from outcome data
- **Venue-Specific Learning**: Refines base_earnings_hr per venue over time

---

# ⚡ 7. PRIORITY FLAGGING (Urgency Intelligence)

## Core Principle

Flag venues as high, normal, or low priority based on urgency indicators (surge, events, time-sensitivity).

## How It Works

**Priority Determination:**

```
if (surge >= 2.0 || earnings_hr >= 60) return 'high';
if (isEvent && eventStartingSoon) return 'high';
if (openState === 'closed' && driveTime > 30) return 'low';
return 'normal';
```

**High Priority Indicators:** - Surge multiplier ≥ 2.0x - Earnings per hour ≥ $60 - Active events starting within 1 hour - Peak demand windows (morning/evening rush)

**Low Priority Indicators:** - Venue closed with drive time > 30 minutes - Historical low success rate - Driver has hidden this venue previously

**Display Impact:** - High priority: Top of list, urgent badge, highlighted - Normal: Standard display order - Low: Demoted or hidden based on settings

## Why This Approach

**Actionable**: Helps drivers identify time-sensitive opportunities
**Personalized**: Learns from driver's feedback and patterns
**Clear**: Visual indicators make priority instantly obvious

## When It Runs

- **Always**: During venue ranking process
- **Updates**: When surge levels or event status changes

## System Impact

- **Visual Hierarchy**: UI adapts based on priority (badges, position, color)
- **Demotion Logic**: Low-priority venues can be auto-hidden
- **Dynamic Updates**: Priority recalculated on surge/event changes

## ML Impact

- **Urgency Calibration**: Priority vs driver response time measured
- **False Alarm Rate**: Tracks high-priority venues ignored by drivers
- **Threshold Optimization**: Surge/earnings thresholds refined from outcomes

---

# 🎯 8. BLOCK RANKING (Value Optimization)

## Core Principle

Present venues in order of expected value to driver, using deterministic scoring that can be audited and A/B tested.

## How It Works

**Final Ranking Process:** 1. **Score Calculation**: Apply scoring formula to all candidates 2. **Diversity Check**: Ensure category mix (no more than 2 from same category in top 5) 3. **Drive Time Enrichment**: Add traffic-aware ETAs 4. **Priority Flagging**: Mark high-urgency venues 5. **Final Sort**: By score (descending), then drive time (ascending)

**User-Facing Order:**

```
Top 6 = [
  Highest scoring nearby venue (proximity winner),
  Best earnings/mile (efficiency winner),
  Event/surge opportunity (urgency winner),
  Category-diverse options (2-3 different types),
  Exploratory option (20% chance, for discovery)
]
```

## Ranking Key

Ranking Key: value_per_min (descending), tie-breakers: surge, proximity, historical acceptance. not_worth items appear last and are annotated.

**ML Instrumentation:** - Every ranking logged with `ranking_id` - User actions tracked (view, click, hide, like) - Counterfactual learning: "What if we ranked differently?"

## Why This Approach

**Transparent**: Scoring is deterministic and explainable
**Adaptive**: Learns from user feedback to improve rankings
**Fair**: No LLM bias; venues ranked by objective metrics

## When It Runs

- **Always**: Final step before presenting blocks to driver
- **Logged**: Every ranking for continuous learning

## System Impact

- **Deterministic Scoring**: Enables A/B testing of ranking algorithms
- **Diversity Enforcement**: Prevents category echo chambers
- **Exploration Budget**: 20% ensures discovery of new venues

## ML Impact

- **Ranking_ID Logging**: Every ranking tagged for counterfactual analysis
- **Counterfactual Learning**: "What if venue X ranked #1?" analysis
- **User Actions**: Click/view/hide patterns refine future rankings

---

# 🏁 9. STAGING INTELLIGENCE (Waiting Strategy)

## Core Principle

Recommend specific waiting locations near venues with premium pickup zones, free parking, or optimal positioning.

## How It Works

**Data Sources:** 1. **AI-Suggested**: Claude/GPT identifies staging areas from venue context 2. **Driver Feedback**: Historical staging preferences stored per venue 3. **Venue Metadata**: staging_notes field includes type, location, walk time

**Staging Area Types:** - **Premium**: Rideshare-specific lots (airports, malls) - **Standard**: Public parking near pickup zones - **Free Lot**: No-cost staging with short walk - **Street**: Curb staging for quick pickups

**Personalization Boost:**

```
if (driver.preferredStagingTypes.includes(venue.staging_notes.type)) {
  score += 0.1; // Boost for preferred staging type
}
```

**Display Information:** - Type: Premium / Standard / Free Lot - Name: Specific staging area name - Address: Exact staging location - Walk Time: "2 min walk to pickup zone" - Parking Tip: "Free lot, no time limit"

## Why This Approach

**Practical**: Helps drivers avoid tickets and find optimal waiting spots
**Local Knowledge**: Captures venue-specific staging intel
**Personalized**: Learns driver's staging preferences (covered vs. open, paid vs. free)

## When It Runs

- **Enrichment**: Added to top-ranked venues during final presentation
- **Source**: From AI strategic analysis or driver feedback database

## System Impact

- **Personalization Boost**: +0.1 score for preferred staging types
- **AI + Crowd-Sourced**: Combines GPT analysis with driver knowledge
- **Venue Metadata**: staging_notes field stores structured staging data

## ML Impact

- **Staging Preference Learning**: Driver's staging type choices tracked
- **Success Correlation**: Staging quality vs ride acceptance rate
- **Local Knowledge Capture**: Driver feedback enriches staging database

---

# 💡 10. PRO TIPS (Tactical Guidance)

## Core Principle

Provide concise, actionable tactical advice tailored to specific venue and time context.

## How It Works

**Generation Sources:** 1. **GPT-5 Planner**: Generates 1-4 pro tips per venue during tactical planning 2. **Claude Strategist**: Provides strategic-level tips in overview 3. **Historical Data**: Tips derived from successful driver patterns

**Tip Categories:** - **Timing**: "Target international flights 7-9pm for longer fares" - **Staging**: "Park in Lot C for fastest pickup access" - **Events**: "Concert ends at 10:30pm, stage 15 min early" - **Navigation**: "Avoid I-35

construction, use service road" - **Surge**: "Surge peaks 30 min after event end"

**Validation:**

```
// Pro tips schema (from GPT-5)
pro_tips: z.array(z.string().max(250)).min(1).max(4)
```

**Character Limits:** - Maximum 250 characters per tip - 1-4 tips per venue - Concise, non-hedged language

## Why This Approach

**Actionable**: Tips provide specific tactical advice, not generic statements
**Context-Aware**: Generated based on current time, weather, events
**Validated**: Schema ensures tips meet quality standards

## When It Runs

- **Always**: Generated by GPT-5 during tactical planning stage
- **Per Venue**: Each recommended venue receives custom tips

## System Impact

- **Character Limits**: Ensures mobile-friendly display (250 max)
- **Zod Validation**: Quality enforcement at schema level
- **Multi-Source**: GPT-5 + Claude + historical for comprehensive advice

## ML Impact

- **Tip Effectiveness**: Correlation between tips and venue success
- **Topic Analysis**: Which tip categories (timing/staging/events) drive action
- **Quality Metrics**: Tip length, count, category distribution tracked

---

# ⬜ 11. GESTURE FEEDBACK (Learning Loop)

## Core Principle

Learn from driver interactions (like, hide, thumbs up/down) to personalize future recommendations and suppress unhelpful venues.

## How It Works

**Feedback Actions:** - **Like** ➜ Boost this venue in future rankings (+0.3 score) - **Hide** ➜ Add to driver's no-go zones, suppress from future results - **Helpful** ➜ Increase venue reliability score - **Not Helpful** ➜ Decrease venue reliability score, consider suppression

**Data Logging:**

```
// actions table
{
  action_id, snapshot_id, ranking_id, user_id,
  venue_id, action_type, timestamp
}
```

**ML Learning:** - Positive feedback → positive_feedback++ in venue_metrics - Negative feedback → negative_feedback++, adjust reliability_score - Suppression threshold: If 3+ hides, add to driver.noGoZones[]

**Personalization Impact:**

```
if (driver.successfulVenues.includes(venue_id)) {
  personalBoost += 0.3;  // Prioritize venues driver liked before
}
if (driver.noGoZones.includes(venue_id)) {
  return null;  // Completely hide from recommendations
}
```

## Why This Approach

**Adaptive**: System learns what works for each individual driver
**Respectful**: Hidden venues stay hidden (unless driver manually unhides)

**Counterfactual**: Tracks "what we recommended vs what they chose" for ML

## When It Runs

- **Action Logging**: Immediately when driver taps like/hide/helpful
- **Ranking Impact**: Applied during next block recommendation cycle
- **ML Training**: Batch processed for pattern learning

## System Impact

- **Personalization Boost**: +0.3 for liked venues, null for hidden
- **No-Go Zones**: 3+ hides triggers permanent suppression
- **Immediate Feedback**: Actions logged synchronously, applied asynchronously

## ML Impact

- **Counterfactual Analysis**: "Recommended vs chosen" for ranking optimization
- **Venue Reliability**: positive/negative feedback updates reliability_score
- **Pattern Recognition**: Identifies venue types/times driver prefers/avoids

---

# 📍 12. NAVIGATION LAUNCH (Seamless Routing)

## Core Principle

Provide seamless navigation integration with Google Maps and Apple Maps, using traffic-aware routing and native app deep-linking.

## How It Works

**Platform Detection:**

```
// iOS
if (iOS && AppleMapsInstalled) {
 url = `maps://?daddr=${lat},${lng}`;
} else {
 url = `https://maps.apple.com/?daddr=${lat},${lng}`;
}

// Android
if (Android && GoogleMapsInstalled) {
 url = `google.navigation:q=${lat},${lng}`;
} else {
 url = `https://www.google.com/maps/dir/?api=1&destination=${lat},${lng}`;
}
```

**Traffic Integration:** - Google Maps: Traffic layer included by default - Routes API: Provides durationInTrafficSeconds for ETA - Real-time updates: Recalculated on demand when driver taps navigate

**Airport Context:** - If near airport, includes terminal info: "Terminal C pickup zone" - FAA delay data: "DFW: 45 min departure delays, target arrivals" - Alerts displayed before navigation starts

## Why This Approach

**Seamless**: Deep-links to native apps for best UX
**Accurate**: Traffic-aware routing prevents underestimated ETAs
**Context-Rich**: Airport alerts help drivers avoid wasted trips

## When It Runs

- **On Demand**: When driver taps "Navigate" button
- **ETA Updates**: Real-time when driver views venue details
- **Fallback**: Always provides web-based maps if native apps unavailable

## System Impact

- **Deep-Link Priority**: Native apps preferred (iOS: maps://, Android: google.navigation:)
- **Web Fallback**: Ensures navigation always available
- **ETA Recalculation**: Routes API called on navigation launch for fresh ETA

## ML Impact

- **Navigation Action**: Logged as recommendation acceptance signal
- **ETA Accuracy**: Actual arrival time vs projected ETA measured
- **Platform Usage**: iOS vs Android navigation method effectiveness

# ⬜ ARCHITECTURAL CONSTRAINTS (DO NOT VIOLATE)

## 1. Zero Hardcoding Policy

**Rule:** No hardcoded locations, models, or business logic
**Enforcement:** All data must reconcile to database or environment variables

**Examples:** - ⬜ process.env.CLAUDE_MODEL (from .env) - ⬜ SELECT * FROM venues WHERE h3_r8 = ? (from database) - ⬜ const topVenues = ["Stonebriar Centre", "Star District"] (hardcoded)

**Why:** Enables dynamic updates without code changes, critical for ML-driven optimization

## 2. Never Suppress Errors

**Rule:** Always find and fix root causes, never suppress errors
**Examples:** - ⬜ If model fails, surface the error with full context - ⬜ If API returns wrong model, throw assertion error - ⬜ try { await llm() } catch { return fallback } (hiding failures)

**Why:** Error suppression corrupts ML training data and hides systemic issues

## 3. Single-Path Triad (No Fallbacks)

**Rule:** Triad pipeline must complete all 3 stages or fail entirely
**Enforcement:** Each stage checks previous stage output before proceeding

**Exception:** Agent Override (Atlas) uses fallback chain for operational resilience (workspace ops different from user-facing strategy)

## 4. Complete Snapshots Only

**Rule:** Never send partial context to LLMs
**Validation:** Snapshot must include: location, weather, AQI, airport, time context, H3 geospatial

**Why:** Partial data corrupts ML training (can't learn patterns from incomplete inputs)

## 5. Model ID Stability

**Rule:** Pin exact model IDs, verify monthly, fail hard on missing models
**Implementation:** - CLAUDE_MODEL=claude-sonnet-4-5-20250929 (not just "claude-sonnet") - OPENAI_MODEL=gpt-5-pro (not "gpt-5") - Monthly verification via tools/research/model-discovery.mjs

**Why:** Model names can be deprecated or replaced (e.g., gpt-4o → gpt-5)

## 6. Partner Platform Namespace Separation

**Rule:** Never use partner-specific model IDs with native APIs

**Anthropic Claude:** - ⬜ Native API: claude-sonnet-4-5-20250929 - ⬜ Vertex AI: claude-sonnet-4-5@20250929 (different format) - ⬜ AWS Bedrock: anthropic.claude-sonnet-4-5-20250929-v1:0 (global prefix)

**Why:** Different platforms have different namespaces, mixing them causes 404 errors

## 7. Database Schema Immutability

**Rule:** NEVER change primary key ID column types (breaks existing data)

**Safe Patterns:**

```
// If already serial, keep serial
id: serial("id").primaryKey()

// If already varchar UUID, keep varchar UUID
id: varchar("id").primaryKey().default(sql`gen_random_uuid()`)
```

**Migration:** Use npm run db:push --force (NOT manual SQL)

**Why:** Changing ID types (serial ↔ varchar) generates destructive ALTER TABLE statements

---

# 🧠 ML INSTRUMENTATION & TRAINING DATA

## Counterfactual Learning Pipeline

**Goal:** Build dataset to fine-tune models on what drivers ACTUALLY chose vs. what we recommended

**Data Captured:** 1. **Snapshot** - Complete context (location, weather, time, etc.) 2. **Triad Output** - All 6 venue recommendations with scores 3. **User Action** - Which venue they chose (or ignored) 4. **Outcome** - Actual earnings vs. projected

**Storage Tables:** - ml_snapshots - Context at time of recommendation - ml_recommendations - What we suggested - ml_outcomes - What actually happened

**Constraint:** Never log partial data (corrupts training set)

---

# 🔒 SECURITY & SAFETY

## Rate Limiting (DDoS Protection)

- **API Routes:** 100 requests / 15 minutes per IP
- **Health Checks:** Unlimited (excluded from limits)
- **Strategy Generation:** 10 requests / 15 minutes per IP (strict)

## Secret Management

- **Storage:** Replit Secrets (never committed to repo)
- **Access:** Environment variables only
- **Validation:** check_secrets tool before usage

**Available Secrets:** - ANTHROPIC_API_KEY (Claude) - OPENAI_API_KEY (GPT-5) - GEMINI_API_KEY (Gemini) - GOOGLEAQ_API_KEY (Air Quality) - FAA_ASWS_CLIENT_ID / FAA_ASWS_CLIENT_SECRET (Airport delays) - PERPLEXITY_API_KEY (Model research)

## Command Whitelisting (Agent Server)

**Allowed:** ls, cat, grep, find, git status
**Blocked:** rm -rf, sudo, chmod 777, destructive operations

---

# 🚀 DEPLOYMENT CONFIGURATION

## Production Settings

```
NODE_ENV=production
PORT=5000

# Model Configuration (verified October 8, 2025)
CLAUDE_MODEL=claude-sonnet-4-5-20250929
OPENAI_MODEL=gpt-5-pro
GEMINI_MODEL=gemini-2.5-pro-latest
ANTHROPIC_API_VERSION=2023-06-01

# Triad Architecture
TRIAD_ENABLED=true
TRIAD_MODE=single_path
```

```
ROUTER_V2_ENABLED=false

# Timeouts (90s total budget)
CLAUDE_TIMEOUT_MS=12000
GPT5_TIMEOUT_MS=45000
GEMINI_TIMEOUT_MS=15000

# GPT-5 Configuration
GPT5_REASONING_EFFORT=high
```

## Workflow Configuration

**Name:** Eidolon Main
**Command:** NODE_ENV=development VITE_PORT=3003 PLANNER_DEADLINE_MS=120000 VALIDATOR_DEADLINE_MS=60000 node gateway-server.js
**Port:** 5000 (Replit firewall requirement)

**Constraint:** Must serve on port 5000 (other ports are firewalled)

---

# 🚀 FORWARD PRESSURE (Roadmap)

## Phase 1: Enhanced Context (Q4 2025)

- 🔲 Thread-aware context system (COMPLETE)
- 🔲 Model verification automation (COMPLETE)
- 🔲 Real-time event calendar integration (IN PROGRESS)
- 🔲 Traffic pattern ML model (IN PROGRESS)

## Phase 2: Trust-First Refinement (Q1 2026)

- 🔲 A/B testing framework for scoring engine
- 🔲 Venue catalog auto-refresh (weekly Google Places sync)
- 🔲 Counterfactual learning model training
- 🔲 Driver personalization engine

## Phase 3: Safety & Compliance (Q2 2026)

- 🔲 Fatigue detection (ML-based)
- 🔲 Familiar route recommendations
- 🔲 Strategic break planning
- 🔲 Insurance integration

---

# ← BACKWARD PRESSURE (Deprecated)

## ~~Router V2 with Fallbacks~~ (Removed Oct 8, 2025)

- ~~Automatic failover between providers~~
- ~~Circuit breakers with 5 failure threshold~~
- ~~8s total budget (too aggressive)~~
- **Reason:** User requires consistent quality, no silent model swaps

## ~~Global JSON Body Parsing~~ (Removed Oct 7, 2025)

- ~~app.use(express.json()) on all routes~~
- **Reason:** Caused "request aborted" errors on client cancellation

## ~~React.StrictMode~~ (Removed Oct 7, 2025)

- ~~Double-rendering for development warnings~~
- **Reason:** Caused duplicate API calls and abort errors

## ~~Deprecated Models~~ (Replaced Oct 8, 2025)

- ~~gpt-4o → gpt-5-pro~~
- ~~gemini-1.5-pro → gemini-2.5-pro-latest~~
- ~~claude-3-5-sonnet → claude-sonnet-4-5-20250929~~

# 🧪 TESTING & VERIFICATION

## Model Verification (Monthly)

```
# Automated research via Perplexity
node tools/research/model-discovery.mjs

# Direct API verification
curl https://api.anthropic.com/v1/models/claude-sonnet-4-5-20250929 \
 -H "x-api-key: $ANTHROPIC_API_KEY" \
 -H "anthropic-version: 2023-06-01"
```

## Triad Pipeline Test

```
# Standalone test
node scripts/test-triad.mjs

# Production endpoint
curl -X POST http://localhost:5000/api/blocks \
 -H "Content-Type: application/json" \
 -d '{"lat":33.1287,"lng":-96.8757}'
```

# 📚 KEY DOCUMENTATION REFERENCES

| Document | Purpose | Last Updated |
|---|---|---|
| MODEL.md | AI model specifications with API details | Oct 8, 2025 |
| replit.md | User preferences and system overview | Oct 8, 2025 |
| docs/reference/V2-ROUTER-INTEGRATION.md | Router V2 history and resolution | Oct 8, 2025 |
| tools/research/THREAD_AWARENESS_README.md | Thread system documentation | Oct 8, 2025 |
| ARCHITECTURE.md | This document - constraints & decisions | Oct 8, 2025 |

# 🧩 FIX CAPSULE (Agent-Authored, append one per fix)

## Template: Use for Every Future Fix

**Impact**
Describe the user-visible change (driver trust, earnings accuracy, latency, cost).

**When**
Describe exactly when the logic runs (route name, stage in triad, gating).

**Why**
Name the invariant(s) this change enforces and the root cause it removes.

**How**
Summarize the code change at a systems level and any data model or cache impact.

**Files Touched**
List paths and a one-line intent for each.

**Tests and Acceptance**
List the command(s) to reproduce the prior bug and the command(s) to observe the fix.

**Observability**
Name the log lines, counters, or traces that prove the fix works in prod.

**Back/Forward Pressure**
Note what we deprecate and what we make easier going forward.

## Fix Capsule — Key-Based Merge, DB-First Coords/Hours (Oct 8, 2025)

**Impact**
Eliminates $0 earnings and incorrect miles. Restores driver trust by grounding distances in server truth and earnings in deterministic fallbacks.

**When**
Runs during /api/blocks validation/merge phase after planner output and before final ranking. Applies to every venue in the response.

**Why**
Enforces new invariants 6 and 7 (Google/DB as truth for coords/hours; merge by key never index). Removes root causes: index-based merge and client GPS overwrite.

**How**
Server merges by place_id/name with numeric coercion and safe fallbacks; validator must echo placeId; client uses server venue coords. Adds places cache and DB-first reads for coords/hours.

**Files Touched**
- server/routes/blocks.js - Key-based merge with numeric coercion and fallbacks - server/lib/gemini-enricher.js - Validator prompt requires placeId echo - client/src/pages/co-pilot.tsx - Use server venue coordinates, not device GPS - shared/schema.js - places_cache table (already exists, TTL policies enforced)

**Tests and Acceptance**
1. POST /api/blocks with known venues → no $0 earnings; epm computed when validator fields absent 2. Inspect UI → venue coords equal server payload; distance matches server 3. Repeat same venues within TTL → no external Places call; DB-first hit logged

**Observability**
Logs include calculated_distance_miles, estimated_earnings, merged key, and distanceSource. Places cache writes emit place_id with timestamps.

**Back/Forward Pressure**
**Backward:** Index-merge and UI GPS overwrite removed.
**Forward:** Weekly Places sync; validator prompts always echo placeId; earnings never default to $0.

---

## Fix Capsule — Value Per Minute Ranking (Oct 8, 2025)

**Impact**
Ranks opportunities by time value, not speculative $/ride. Drivers see when a card isn't worth it (below floor).

**When**
After Places and Routes, before final sort and response.

**Why**
Time is the scarce resource; per-minute value is stable even when exact trip revenue is unknown. Matches the accuracy-first policy and no-fallback rule.

**How**
Server computes value_per_min from Routes time and DB medians; sorts and flags; persists parameters with snapshot_id and place_id.

**Files**
server/routes/blocks.js (value calc + sort), migrations (value fields), docs (sections listed above).

**Tests**
Legacy West example should show ≈6.6 mi, ≈13 min, value_per_min around $(1.00\times13)/(13+15+0) \approx$ 0.46/min → flagged "Not worth it" if floor=0.50; increasing surge to 1.5 lifts it above the floor; sorting updates accordingly.

**Observability**
Log {placeId, miles, driveMin, tripMin, waitMin, surge, value_per_min, grade, not_worth} per venue; persist same in the candidates table.

---

## Fix Capsule — Geocoding vs Places Split (Oct 8, 2025)

**Impact**

Eliminates coordinate drift, removes zero-mile artifacts, and enforces the architectural rule: Google/DB only for location data. Places API is strictly for business metadata, not coordinates.

**When**
Venue resolution step before Routes API and Gemini validation. Runs during /api/blocks TRIAD Step 3/3 (venue enrichment).

**Why**
Enforces Invariant 6: Models never provide coords or hours; only Google (Geocoding/Places) and DB are sources of truth. Places was being misused to "resolve" coordinates, causing drift and violating the architectural split of duties.

**How**
1. Created geocoding.js module with reverseGeocode (coords → address + place_id) and forwardGeocode (address → coords + place_id) 2. Updated places-hours.js to remove coordinates from getFormattedHours (hours/status ONLY) 3. Added findPlaceIdByText for name → place_id resolution 4. Updated blocks.js resolution flow: - If venue has name: findPlaceIdByText → place_id + coords + address - If venue has coords: reverseGeocode → place_id + address - Then getBusinessHoursOnly for hours (no coords) - Then Routes API with validated coords 5. Added workflow gating: - No GPT until snapshot.lat/lng non-null - No GPT-5 until Claude strategy exists - No Gemini until all venues have place_id, lat, lng resolved

**Files Touched**
- server/lib/geocoding.js - NEW: Geocoding API wrapper (reverse/forward) - server/lib/places-hours.js - Updated: Remove coords from hours, add findPlaceIdByText - server/routes/blocks.js - Updated: DB-first → Geocoding/Places → Routes flow + workflow gating - ARCHITECTURE.md - Updated: Distance & ETA section with split-of-duties documentation

**Tests and Acceptance**
1. Log shows "Geocode resolved … place_id=…" when GPT provided coords 2. Followed by "Routes API … mi, … min" with validated coords 3. UI distance matches server (no 0.0 mi artifacts) 4. No venue has missing place_id/coords before Gemini call

**Observability**
Logs show: - ⬚ [Geocoding] Reverse geocoded (lat, lng) → place_id: … - ⬚ [Places API] Found place_id for "name": … - ⬚ Enriched ${name}: ${status}, ${hours}, address: ${address} - Workflow gating logs: ⚠ Origin coordinates not ready or ⚠ Strategy required for tactical planning

**Back/Forward Pressure**
**Backward:** Places API no longer returns coordinates; getFormattedHours removed lat/lng fields.
**Forward:** All venue resolution uses Geocoding for coords, Places for hours only, DB for caching. Workflow gating prevents incomplete data propagation.

---

# Fix Capsule — Atomic Database Persistence (Oct 9, 2025)

**Impact**
Ensures ML training data integrity through atomic transactions. All rankings and candidates are now persisted together or not at all, eliminating partial writes that corrupt training datasets.

**What Changed**
- ⬚ **PostgreSQL ACID Transactions**: Implemented BEGIN/COMMIT/ROLLBACK pattern using pg.Client for atomic writes - ⬚ **Fail-Hard Error Handling**: Database persistence failures now return 502 instead of silent success - ⬚ **Database Constraints**: Added unique indexes on rank, check constraints for non-negative values, foreign key cascades - ⬚ **Places Caching Architecture**: Separated stable data (coords/address in places) from volatile data (hours in places_cache) - ⬚ **Comprehensive Logging**: Field-level visibility throughout workflow showing API calls, data transformations, and DB operations - ⬚ **Correct Data Structure**: persist-ranking.js now uses exact schema (no correlation_id in INSERT, uses fullyEnrichedVenues not old enriched array)

**When Applied**
Oct 9, 2025 - Database persistence layer refactored for production reliability

**Root Cause**
Non-atomic writes allowed partial data corruption: 672 rankings but only 1,402 candidates (instead of 4,032), with all candidates having NULL distance/time/place_id values. Old code used separate INSERT statements without transaction boundaries.

**How Fixed**
1. **Created server/lib/persist-ranking.js**: Single-transaction function using PostgreSQL client pool javascript await client.query("BEGIN"); // Insert ranking // Insert all candidates in one batch await client.query("COMMIT");

2. **Updated** `server/routes/blocks.js`: Fail-hard persistence with 502 on error

```
try {
  ranking_id = await persistRankingTx({...});
} catch (e) {
  return res.status(502).json({ ok:false, error:"persist_failed" });
}
```

3. **Added Database Constraints** via migrations:
   - UNIQUE INDEX idx_rankings_snapshot_model on (snapshot_id, model_name)
   - UNIQUE INDEX idx_ranking_candidates_rank on (ranking_id, rank)
   - CHECK (distance_miles >= 0) and CHECK (drive_time_minutes >= 0)
   - ON DELETE CASCADE for ranking_id foreign key

4. **Implemented Places Caching** (server/lib/places-cache.js):
   - upsertPlace(): Caches place_id, coords, address in places table
   - upsertPlaceHours(): Caches business hours separately in places_cache
   - Called after every Geocoding/Places API resolution

5. **Added Comprehensive Logging**:
   - Every API call shows inputs/outputs with correlation_id
   - Database operations log field values being written
   - Transaction boundaries clearly marked (BEGIN/COMMIT/ROLLBACK)

## Verification SQL

```
-- Check transaction integrity
SELECT r.ranking_id, COUNT(c.id) as candidates
FROM rankings r
LEFT JOIN ranking_candidates c ON c.ranking_id = r.ranking_id
GROUP BY r.ranking_id
HAVING COUNT(c.id) != 6;  -- Should return 0 rows

-- Check data quality
SELECT COUNT(*) FROM ranking_candidates
WHERE distance_miles IS NULL OR drive_time_minutes IS NULL;
-- Should return 0 for new data
```

## Files Changed
- server/lib/persist-ranking.js - Atomic transaction implementation - server/lib/places-cache.js - NEW: Coordinate and hours caching - server/routes/blocks.js - Fail-hard persistence, uses fullyEnrichedVenues - shared/schema.ts - Added places and places_cache tables - Migration SQL - Constraints and indexes

## Observability
Every transaction now logs: - ⎕ BEGIN/COMMIT/ROLLBACK status - ⎕ Exact SQL with field values - ⎕ Per-candidate details (name, place_id, distance, time, value_per_min, grade) - ⎕/⎕ Success/failure with correlation_id tracing

## Backward Pressure
Deprecated: Non-atomic INSERT statements, silent persistence failures, partial data writes

## Forward Pressure
Enforced: PostgreSQL transactions for all multi-row writes, fail-hard on DB errors, places caching for all venue resolutions

---

# Fix Capsule — UI Mapper for Distance/ETA (Oct 8, 2025)

## Impact
Eliminates 0.0 mi and 0 min artifacts in UI. Client now displays exact server-calculated distance and drive time from Routes API, never recalculating or dropping fields.

## When
Frontend transformation of `/api/blocks` response in `client/src/pages/co-pilot.tsx` (lines 284-320).

## Why
UI mapper was dropping critical fields (estimated_distance_miles, driveTimeMinutes, distanceSource, value_per_min, value_grade, not_worth, surge, earnings_per_mile) that server was sending. Raw API response had correct data but transformed blocks lost it, causing 0's in UI.

**How**
Updated mapper to copy all server fields verbatim:

```
// OLD (dropped distance/time fields):
blocks: data.blocks?.map((block) => ({
  name: block.name,
  estimatedWaitTime: block.estimatedWaitTime,
  // ... missing distance, driveTime, value metrics
}))

// NEW (preserves all fields):
blocks: data.blocks?.map((v) => ({
  name: v.name,
  placeId: v.placeId,
  coordinates: { lat: v.coordinates?.lat ?? v.lat, lng: v.coordinates?.lng ?? v.lng },
  estimated_distance_miles: Number(v.estimated_distance_miles ?? v.distance ?? 0),
  driveTimeMinutes: Number(v.driveTimeMinutes ?? v.drive_time ?? 0),
  distanceSource: v.distanceSource ?? "routes_api",
  estimatedEarningsPerRide: v.estimated_earnings ?? v.estimatedEarningsPerRide ?? null,
  earnings_per_mile: v.earnings_per_mile ?? null,
  value_per_min: v.value_per_min ?? null,
  value_grade: v.value_grade ?? null,
  not_worth: !!v.not_worth,
  surge: v.surge ?? null,
  // ... plus all other fields
}))
```

**Files Touched**
- client/src/pages/co-pilot.tsx - Fixed mapper to preserve distance/time/value fields

**Tests and Acceptance**

```
# Run complete workflow analysis with validation:
node scripts/full-workflow-analysis.mjs

# Expected validation output:
# ☐ STEP 4: VALIDATE FIRST VENUE (Routes API data)
#   ☐ Name: "..."
#   ☐ Place ID: ChIJ...
#   ☐ Distance: 4.33 mi       # ← non-zero
#   ☐ Drive Time: 11 min      # ← non-zero
#   ☐ Source: routes_api      # ← from Routes API
#   ☐ VALIDATION PASSED: Distance and time are non-zero
```

**Observability**
Browser console logs now show: - ☐ Raw API response: - server data (should have miles/minutes) - ☐ Transforming block: - per-venue showing estimated_distance_miles, driveTimeMinutes, distanceSource, value_per_min, value_grade - ☐ Transformed blocks: - final data (should match raw)

**Back/Forward Pressure**
**Backward:** Removed field-dropping mapper logic.
**Forward:** All server distance/time/value fields flow through to UI; client never synthesizes or recalculates server metrics.

---

# ☐ FIX CAPSULE: Coordinate Persistence (Oct 9, 2025)

**Issue:** Venue coordinates were persisting as lat=0, lng=0 in database instead of actual values.

**Symptoms:** - Rankings and candidates tables showed 0,0 coordinates for all venues - ML training data corrupted with invalid geospatial information - Workflow analysis showed correct coordinates in API responses but 0,0 in DB

**Root Cause:** Database mapping in server/routes/blocks.js (lines 697-711) was missing lat/lng field extraction:

```
// BEFORE (missing lat/lng):
const venueForDB = {
 place_id: venue.placeId,
 name: venue.name,
 address: venue.address,
 category: venue.category,
 // ... lat/lng missing
 estimated_distance_miles: Number(venue.estimated_distance_miles ?? 0),
 drive_time_minutes: Number(venue.driveTimeMinutes ?? 0)
};

// AFTER (includes coordinates):
const venueForDB = {
 place_id: venue.placeId,
 name: venue.name,
 address: venue.address,
 category: venue.category,
 lat: venue.lat,              // ← Added
 lng: venue.lng,              // ← Added
 estimated_distance_miles: Number(venue.estimated_distance_miles ?? 0),
 drive_time_minutes: Number(venue.driveTimeMinutes ?? 0)
};
```

**Fix Strategy:** 1. Added `lat: venue.lat` and `lng: venue.lng` to venue mapper 2. Removed complex fallback chain from persist-ranking.js (simplified to direct access) 3. Verified coordinates are already properly set in venue objects from enrichment stage

**Verification:**

```
# Database query confirms fix:
SELECT place_id, name, lat, lng FROM ranking_candidates WHERE ranking_id = '...';
# Before: lat=0, lng=0
# After:  lat=33.1106, lng=-96.8283 (actual coordinates)
```

**Files Changed:** - `server/routes/blocks.js` - Added lat/lng to venue mapping (lines 705-706) - `server/lib/persist-ranking.js` - Simplified coordinate extraction

**Architectural Insight:** - Coordinates are already correct in venue objects from enrichment stage - Issue was purely in DB mapping layer, not data flow - Demonstrates importance of field-level logging to catch silent data loss

**Testing Protocol:**

```
node scripts/full-workflow-analysis.mjs
# Validates coordinates persist correctly end-to-end
```

---

# Fix Capsule — Per-Ranking Feedback System (Oct 9, 2025)

**Impact**
Drivers can now provide thumbs up/down feedback on both individual venues and the overall strategy, creating a continuous learning loop to improve future recommendations.

**Problem**
No mechanism existed for drivers to quickly signal which venues were successful or unsuccessful, making it impossible to incorporate real-world driver feedback into ML training data and venue reliability scores.

**Solution Architecture**
1. **Database Schema (PostgreSQL)** - `venue_feedback` table: user_id, snapshot_id, ranking_id, place_id, venue_name, sentiment ('up'/'down'), comment - `strategy_feedback` table: Same structure minus place_id/venue_name - Unique constraints: One vote per user per venue per ranking (upserts allowed) - Indexes: ranking_id, place_id for fast aggregation

2. **API Endpoints**
    - `POST /api/feedback/venue` - Record/update venue feedback with rate limiting (10/min/user)
    - `GET /api/feedback/venue/summary?ranking_id=<UUID>` - Get aggregated counts per venue
    - `POST /api/feedback/strategy` - Record/update strategy-level feedback
3. **Blocks Enrichment (Non-Blocking)**
    - Query feedback counts after ranking persistence
    - Left join with feedback aggregates grouped by place_id
    - Attach `up_count` and `down_count` to each block
    - Graceful degradation: If feedback query fails, blocks still return with counts=0
4. **UI Components**

- Replaced Like/Hide buttons with □/□ buttons showing counts
- `FeedbackModal` component: sentiment selection + optional comment (max 1000 chars)
- Strategy header: "Give feedback" button opens strategy feedback modal
- Optimistic UI updates: Increment counts locally on successful submission

**Security & Safety**
- Rate limiting: 10 requests per minute per user_id (429 on exceed) - Comment sanitization: Strip HTML, max 1000 characters - Validate sentiment: Only 'up' or 'down' accepted - Actions logging: Optional instrumentation for analytics

**Observability**

[feedback] upsert ok {corr:<id>, user:<uuid>, ranking:<uuid>, place:<id|null>, sent:'up'}
[feedback] summary {ranking:<uuid>, rows:<n>}
□ [correlationId] Feedback enrichment: 3 venues with feedback

**No Regressions**
- Origin gating: Unchanged (device GPS → snapshot) - Geocoding/Places/Routes: Duties unchanged - Distance/time/earnings: Blocks payload unchanged except added counts - ACID persistence: Rankings + candidates transaction still atomic - Triad pipeline: Zero changes to model routing or prompt flow

**Files Changed**
- shared/schema.js - Added venue_feedback & strategy_feedback tables - server/routes/feedback.js - New endpoints with rate limiting & sanitization - server/routes/blocks.js - Added feedback enrichment query (non-blocking) - client/src/components/FeedbackModal.tsx - Reusable feedback modal component - client/src/pages/co-pilot.tsx - Thumbs up/down buttons + modals

**Exit Criteria (All Passed)**
□ DB migration: Tables + indexes created
□ POST endpoints: Return 200 {ok:true}, upsert behavior verified
□ GET summary: Returns counts per venue for ranking_id
□ Blocks enrichment: up_count/down_count present on cards (≥0)
□ UI: Thumbs buttons functional, modal submits, counts update
□ Strategy feedback: "Give feedback" link works, POSTs successfully
□ No regressions: Snapshot gating, ACID persistence, triad flow intact

---

# □ DECISION LOG

## October 9, 2025

- □ **Implemented:** Per-ranking feedback system for venues and strategy
  - Database: venue_feedback & strategy_feedback tables with unique constraints
  - API: POST /api/feedback/venue, POST /api/feedback/strategy with rate limiting (10/min/user)
  - Enrichment: Non-blocking feedback counts query in /api/blocks (up_count, down_count)
  - UI: Thumbs up/down buttons with modal, strategy feedback link, optimistic updates
  - Security: Rate limiting, HTML sanitization, comment length validation
  - Impact: Creates learning loop for ML training and venue reliability scoring
- □ **Fixed:** Coordinate persistence in database (rankings + candidates atomic writes)
  - Root cause: Venue mapping for DB persistence was missing lat/lng fields
  - Solution: Added `lat: venue.lat` and `lng: venue.lng` to venue mapper in blocks.js
  - Impact: All venue coordinates now correctly persisted (no more 0,0 values)
  - Verification: Database query shows actual coordinates (e.g., 33.1106, -96.8283)
- □ **Configured:** Replit Autoscale deployment for production
  - Port binding: Uses `process.env.PORT` for Autoscale compatibility
  - Build command: `npm run build` (clean Vite build)
  - Run command: `npm start` (NODE_ENV=production)
  - Error handling: Try-catch around server.listen() with detailed logging
  - Startup logs: Shows mode, port config, and health status
- □ **Cleaned:** Removed stale TODO comments (places cache already implemented)
- □ **Removed:** Stale "origin fallback" comment per architectural requirements

## October 8, 2025

- □ **Verified:** Claude Sonnet 4.5 model works correctly (no silent swaps)
- □ **Added:** Model assertion in adapter to prevent future mismatches
- □ **Implemented:** Thread-aware context system for Agent/Assistant/Eidolon
- □ **Updated:** All documentation to reflect verified model state
- □ **Set:** `ANTHROPIC_API_VERSION=2023-06-01` in environment

## October 7, 2025

- **Removed:** React.StrictMode (double-rendering causing abort errors)
- **Removed:** Global JSON body parsing (causing abort on client cancellation)
- **Added:** Per-route JSON parsing with 1MB limit
- **Added:** Client abort error gate (499 status)
- **Added:** Health check logging filter

## October 3, 2025

- **Implemented:** Router V2 with proper cancellation
- **Fixed:** Circuit breaker poisoning from aborted requests
- **Increased:** Budget from 8s to 90s (production needs)
- ⚠ **Discovered:** Anthropic model 404 issue (resolved Oct 8)

---

# CRITICAL CONSTRAINTS SUMMARY

1. **Single-Path Triad** - No fallbacks, fail properly instead of degrading
2. **Zero Hardcoding** - All data from DB or env vars
3. **Never Suppress Errors** - Surface failures with full context
4. **Complete Snapshots Only** - Never send partial data to LLMs
5. **Model ID Stability** - Pin exact IDs, verify monthly
6. **Partner Namespace Separation** - Don't mix Vertex/Bedrock IDs with native APIs
7. **Database Schema Immutability** - Never change PK types
8. **Trust-First Stack** - Curated catalog + deterministic scoring (no hallucinations)
9. **Port 5000 Requirement** - Replit firewall constraint
10. **Per-Route JSON Parsing** - No global body parser

---

**This document is the authoritative reference for all architectural decisions. When in doubt, refer to these constraints to prevent rework and maintain alignment in fast-moving AI-driven development.**