



REMRIN.AI: TECHNICAL ARCHITECTURE & SPECIFICATIONS

Proprietary Technology Overview | v1.0

1. EXECUTIVE TECHNICAL SUMMARY

Remrin.ai is a **High-Utility AI Companion Platform** built on a headless, event-driven architecture. Unlike standard chatbot wrappers, Remrin utilizes a proprietary middleware layer—the **R.E.M. Engine** (Resonant Emotional Memory)—to manage long-term memory retention, dynamic persona injection, and multi-model routing.

Core Stack:

- **Frontend:** HTML5/JS (Lightweight Client), expanding to React Native.
 - **Backend Runtime:** Node.js (Edge Functions).
 - **Database:** Supabase (PostgreSQL) with pgvector extension.
 - **AI Orchestration:** Dynamic Router (DeepSeek V3 / Claude 3.5 Sonnet).
 - **Voice Synthesis:** ElevenLabs Enterprise API (Streaming WebSocket).
-

2. THE R.E.M. ENGINE (Resonant Emotional Memory)

The core IP of Remrin is the **R.E.M. Engine**, a custom memory retrieval pipeline designed to solve the "Catastrophic Forgetting" problem inherent in LLMs.

A. The Vector-Relational Hybrid Search

Standard RAG (Retrieval-Augmented Generation) often fails due to "fuzzy" context. We utilize a **Hybrid Search Strategy**:

1. **Semantic Search:** User input is embedded (using text-embedding-3-small) and queried against the memories vector store for emotional context.
2. **Relational Filtering:** Results are filtered by user_id, persona_id, and time_decay_weight to ensure relevance.
3. **The Re-Ranking Layer:** Retrieved fragments are passed through a lightweight Cross-Encoder to re-rank them by *true relevance* before injection.

B. Dynamic Context Injection

The engine dynamically constructs the System Prompt at runtime (\$T=0\$):

- [Static_Constitution] (The "SoulQR" - Immutable rules).
- [Dynamic_Context] (Top-5 Re-ranked Memories).
- [User_State] (Current emotional valence/history).

3. THE SOULFORGE (Persona Synthesis Protocol)

Remrin utilizes a proprietary prompt-engineering methodology known as the "**Lilly Method**" to generate high-fidelity personas without manual creative writing.

- **Archetype Blending:** The system decomposes a user request (e.g., "Big Blue Bear") into constituent archetypes (e.g., Pooh_0.4 + Baloo_0.3 + Sulley_0.3).
- **Vector Synthesis:** The system retrieves behavioral traits from these archetypes and fuses them into a cohesive **System Instruction Block**.
- **Drift Prevention:** We utilize **Separation Tokens** and **Identity Anchoring** within the prompt structure to prevent the LLM from breaking character during long context windows.

4. INFRASTRUCTURE & SCALABILITY

The platform is designed for horizontal scalability using a **Serverless/Edge** architecture.

- **Stateless API:** All chat interactions are stateless; context is re-hydrated from Supabase on every turn, allowing infinite horizontal scaling of compute nodes.
- **Real-Time State:** We utilize **Supabase Realtime** (WebSockets) to push state changes (e.g., "Aether" balance updates, "SoulSpark" unlocks) to the client instantly without polling.
- **Security (RLS):** Row Level Security is enforced at the database tier. A user can strictly only query vectors associated with their auth.uid().

5. FUTURE ROADMAP (Q3 2025)

- **Local-First Inference:** Exploring WebGPU implementation to run "Lite" models (like Llama-3-8B) directly in the user's browser to reduce cloud inference costs ("Aether") to near-zero.
- **Multi-Modal Vision:** Giving Souls "sight" via image-to-text pipelines.

James Gray | Founder & CEO *Alexandria, Egypt*

gr4y74@gmail.com

+1 216.307.6363

