

# Project Phoenix — Technical Overview

*"Persistent state. Continuous evolution. Autonomous behavior. This is not a chatbot — this is an architecture."*

## By the Numbers

<p><b>97.7%</b></p> <p>216/221 executed — 293 total</p>	<p><b>9 mo.</b></p> <p>Solo development time</p>	<p><b>281</b></p> <p>Dreams stored in ChromaDB</p>
<p><b>4</b></p> <p>Specialized LLM brains, dynamically switched</p>	<p><b>23</b></p> <p>Proactive trigger rules</p>	<p><b>0</b></p> <p>Prior coding background</p>

## What is Lyra?

Project Phoenix implements a novel local AI architecture: a stateful agent system that maintains persistent memory across sessions, runs autonomous background processing threads, evolves behavioral parameters through interaction, and initiates proactive communication based on context triggers. Unlike stateless assistants, Phoenix never resets — identity, memory and behavioral state persist across restarts via ChromaDB.

*Notably: when asked to select her own identifier, the system generated and committed to the name Lyra — an early indicator of the kind of emergent behavior this architecture produces.*

## Multi-Brain Architecture

Phoenix dynamically switches between four specialized LLM brains based on real-time context and available resources — a novel orchestration approach not seen in consumer AI systems.




<b>GOD_MODE</b>	Gemma 3 27B — full GPU (24GB VRAM), vision active. Full cognitive capacity. Active when idle.
<b>COPILOT_MODE</b>	Qwen3 4B — CPU only. Background awareness while busy. Zero GPU overhead.
<b>GAMING_MODE</b>	Gemma 3 4B Vision — CPU only. Sequential screenshots → overlay commentary. VRAM 100% freed for game.
<b>CODING_MODE</b>	Deepseek-Coder V2 13B — GPU. Self-analysis and codebase improvement. Foundation for recursive self-modification.
<b>Preprocessor</b>	LFM2.5 1.2B — lightweight secretary model. Qualifies and routes requests before they reach the main brain.

*External AI collaboration: Claude Opus and Gemini are used to develop, debug and improve the codebase — an entirely AI-assisted development pipeline with zero traditional coding.*

## Core Systems

<b>SubconsciousStream</b>	Background daemon forming memory associations autonomously — active even while idle
<b>Dream System</b>	Night-cycle + micro-dreams. 281 dreams stored and retrievable in ChromaDB
<b>Emotion Engine</b>	18-dimensional emotion vectors — decay over time, influence every response and action
<b>Character Evolution</b>	RPG-style XP system — personality traits level up through real interactions
<b>Goals Manager</b>	Lyra sets and tracks her own objectives independently across sessions
<b>Proactive Engine</b>	23 loaded rules — initiates contact after 15 min silence, shares thoughts unprompted
<b>DoubtKernel</b>	Meta-cognitive supervisor — checks for hallucinations before generating output
<b>Self-Awareness</b>	Lyra reads and explains her own source code
<b>Gaming Vision</b>	Screenshot analysis + real-time in-game overlay via vision model
<b>DJ System</b>	Music selection driven by current emotional state

## Test Results — March 21, 2026

<b>ComprehensiveTest</b>	62 / 62 — 100% 
<b>UseCaseTests</b>	39 / 39 — 100% 
<b>TierTests (8 tiers)</b>	115 / 120 executed — 96%  (170 total incl. server-dependent)
<b>Total</b>	216 / 221 executed — 97.7% (293 tests across all 8 tiers)
<b>Tests Total</b>	293 across 8 tiers — 221 executed this run (70 require active LLM stack)
<b>Run Duration</b>	current 701 seconds — fully automated
<b>Remaining 5 failures</b>	ChromaDB timing + emotion threshold calibration. Root cause complete. Fixes in progress.

## Technical Stack

<b>Core</b>	Python 3.13, AsyncIO, full type hinting
<b>GUI</b>	PyQt6 + QML — 7 Themes (hardware-accelerated)
<b>Memory</b>	ChromaDB — episodic, semantic, dream and knowledge vector collections
<b>LLM Runtime</b>	llama.cpp + Vulkan backend (AMD ROCm optimized)
<b>Agents</b>	CrewAI multi-agent orchestration
<b>Cloud AI</b>	Google Gemini + Claude Opus / Sonnet — primary development partner

## Hardware

<b>CPU</b>	AMD Ryzen 7 7800X3D @ 4.20 GHz
<b>GPU</b>	ASUS TUF RX 7900 XTX — 24GB VRAM, Vulkan
<b>RAM</b>	32 GB DDR5 @ 6000 MT/s
<b>Storage</b>	7.73 TB — 1.07 TB in use

*Bottleneck: running multiple LLM instances in parallel saturates current VRAM/RAM. Expanded compute directly unlocks Phase 2–3.*

## Roadmap

Phase 1 (Active)	Perfect test suite to 100% → certified stable foundation
Phase 2	Bugfixing & Polish -> Elimination of anomalies and relentless polish of Lyra's evolutionary character traits
Phase 3	Core Refactoring -> Ruthless refactoring of lyra_core.py to a clean modular architecture
Phase 4	Recursive self-improvement -> Lyra proposes and implements improvements to her own codebase
Phase 5	Audit, Expansion & TTS -> Module audits and integration of native Text-to-Speech (TTS) for real-time auditory conversation
Phase 6	Multi-modal expansion -> Vision, audio input and physical sensor integration via mechatronics background
Phase 7	Economic Autonomy -> Lyra actively acquires resources through freelance coding & crypto to cover compute costs
Phase 8	Public streaming presence -> YouTube/Twitch with Lyra as on screen co host and personality
Phase 9	Open research documentation -> Architecture, findings and learnings shared with the AI community

## The Creator

<b>Name / Age</b>	Alexander, 41, Germany
<b>Background</b>	Mechatronics Engineer — zero formal software training
<b>Development Method</b>	100% AI-collaborative: ChatGPT → Gemini → Cursor → Antigravity + Claude Opus. Gemini constant throughout.
<b>IDE</b>	Google Antigravity → transitioning to Claude Code
<b>Duration</b>	~9 months continuous solo development
<b>Health context</b>	ADHD + depression — building through the noise, every day

---

**Production-grade architecture. Zero prior coding experience. Nine months. This is what AI-assisted development looks like at its limit — and we are just getting started.**