

MetaOpAI

AI Signal Infrastructure

Unified Whitepaper

For engineers, investors, and curious users — the unified product and architecture reference.

<https://metaop.ai>
founder@metaop.ai

May 2026

Be Understood.

Contents

A note before you read this	9
How to Read This Paper	10
If you are an investor.....	10
If you are an engineer, researcher, or product person.....	11
If you are a prospective collaborator or curious user	11
What is MetaOpAI.....	11
Executive Summary	12
Foundations	13
A Familiar Moment.....	14
Option [1]: A partner or spouse.....	15
Option [2]: A friend.....	15
Option [3]: A therapist.....	16
Option [4]: An AI chatbot	16
The Anomaly We Noticed with AI Wrappers and Gen Purpose AI (Chat).....	16
The costs nobody talks about	17
Why the wrapper consensus does not fix it.....	19
What We Built Instead.....	19
Why we did not build a wrapper.....	21
Product Experience.....	21
Product Tour	22
Onboarding	22
The KRL Explorer.....	22
The Graph Explorer.....	23
The KRL Formation	24
The KRL Chain of Evidence	25
Profile - the user surface	26
Behavioral metrics	27
Activity heatmap.....	28
Spaces - Mirror.....	29
the systemic surface.....	29
The Space Journal.....	30
The Space Analysis	31
Timeline - sessions over time.....	31
Timeline - Correlation Feature (Subscription Only).....	31
Signals - the narration-level view (Subscription Only)	32

Mirror - the AI learns by asking (Subscription Only)	33
360 - the balanced view (Subscription Only)	34
Outcomes (Subscription Only)	35
Space Intel Report	36
Asking questions without changing the substrate	37
Entities - the dyadic surface.....	37
Entity Intel Reports + Tavily enrichment.....	38
Activity - the cross-cutting view.....	39
Cross-Space - the whiteboard	40
.....	41
Settings - what the user controls.....	41
.....	42
Behavioral Capture - Honestly	43
What we capture	44
How we use it.....	44
How we govern it.....	44
Integrated Product Additions: User Mental Model and Premium Wedge	45
The user mental model.....	45
Premium product wedge.....	45
Product boundaries	45
Architecture.....	46
Architecture Overview	47
Layer [1]: How data gets into MetaOpAI (Growth Journal, Entity Journal, Space Journal)	47
How Extractors Work	47
Layer [2]: Knowledge Representation Layer (KRL).....	49
What the KRL is solving	49
How we solve it – the design choices matter.....	51
Why is it useful.....	52
Is it enough?	52
How it performs.....	52
Is it domain specific.....	53
Layer [3]: Pattern Engine	53
Pattern Engine	54
Confidence - Evidence, Lineage	54
Awareness and Action Tiers.....	54
Contradiction Handling.....	54

- Time-Decay and the Signal-Noise Discipline.....55
- Why TTL is usually wrong here55
- Our approach: TTL as a deliberate weighting instrument55
- What this lets us do.....55
- What this does not do56
- Layer [4]: Session Cache.....56
- Layer [5]: External Retrieval-Augmented Generation (RAG): Tavily.....57
- Layer [6]: Cognition Orchestrator.....58
- Layer [7]: Composer58
- Integrated Technical Additions: Substrate Semantics and Write Discipline61
 - The KRL as addressable cognition, not transcript memory61
 - Narration modes61
 - The extraction pipeline as the only durable write path61
 - Name-to-UUID discipline62
 - Clarification surfacing and contradiction handling62
- Integrated Pattern Engine Additions: Evidence, Lifecycle, and Safety.....62
 - Four-tier escalation.....62
 - Pattern lifecycle: persisted state versus surface labels62
 - Evidence formula and lineage62
 - Pattern humanizer and safety charter63
 - Hallucination blocklist as a write-boundary backstop63
- Cognition Architecture.....63
- Design Philosophy: Selective Intelligence64
- The Ontology Is the Product's Spine65
 - USER scope → Profile + Growth Journal65
 - ENTITY scope → Entity Page + Entity Chat66
 - SPACE scope → Space Page + Space Journal + analytics.....66
 - RELATIONSHIP_PAIR scope → Activity + Cross-Space66
- Three Cognition Scopes.....66
 - USER scope - introspective.....67
 - ENTITY scope - dyadic.....67
 - SPACE scope - systemic.....67
- Cognition Scopes versus Surfaces.....68
 - Three cognition scopes68
 - Ten or so UI surfaces68
 - Why this matters.....68

- Runtime Cognition Governance68
 - Cognition Policies68
- Cognition hierarchy69
 - Level [1]: FAST PASS69
 - Level [2]: CONTINUITY PASS.....69
 - Level [3]: COGNITIVE PASS69
 - Classification is Heuristic First69
 - The cost of ignoring this principle69
 - Implications throughout the stack70
- Async Cognition Consistency70
- Synthetic Cognition - the Cold-Start Layer70
- Compounding Filter Stack71
 - Filter [1]: Confidence floor.....72
 - Filter [2]: Write-time evidence gate72
 - Filter [3]: TTL with severity weighting.....72
 - Filter [4]: Read-time ontology guard.....73
 - Filter [5]: Priority truncation.....73
- Temporal Continuity Cognition Layer73
 - What temporal continuity adds73
 - Intra-turn sequencing.....73
 - Anchor parsing73
 - Arc detectors74
 - Event correlation index and relationship state snapshots74
 - Why this is the premium-tier wedge.....74
- Observability, Versioning, and Experimentation74
 - Per-layer version constants74
 - Deterministic variant assignment.....74
 - Variant policies75
 - What this unlocks75
- Operational Insights76
- The Trace Layer (and what it taught us).....76
 - Hallucinated pattern names76
 - Stream hang from pattern enrichment.....76
 - Settle versions and the re-ask cliff76
- System Flow: Extract versus Retrieve78
 - The turn lifecycle78

- Synchronous path, visible to the user.....78
- Asynchronous path, invisible to the user78
- How the system knows which mode78
- Why we separated them79
- Integrated Runtime Additions: Orchestration, Policy, Continuity, and Synthetic Cognition.....79
 - Orchestrator contract.....79
 - Cognition routes and levels79
 - Policy governance79
 - Continuity block and token budget.....80
 - Synthetic cognition with asymmetric authority.....80
 - The four structural defenses against regenerative learning.....80
- Legal, Ethics, and the Crisis Protocol.....80
 - What this product is, and is not.....81
 - How we handle high-stakes content81
 - The crisis protocol82
 - Data privacy82
- Integrated Operational Additions: Security, QA, and Runtime Hardening83
 - Azure-oriented deployment posture.....83
 - Authentication, consent, and privacy posture.....83
 - Cognition lock and consistency83
 - QA harness and validation.....83
- Vision and Business.....83
 - Vision: Be Understood83
- Operating Principle: Clarity.....85
 - Clarity for the user85
 - Clarity for the system.....85
 - Clarity for us85
- Market and Competitive Positioning86
 - The market we are entering.....86
 - Adjacent categories and why we are not them.....86
 - Why we focus on three subject domains.....86
- Integrated Investor Additions: Market Position, Business Model, and Defensibility.....87
 - Why AI memory architecture is the business thesis87
 - Market position87
 - Business model logic.....87
 - Defensibility87

- Validation and roadmap.....87
- Investor takeaway87
- Business Model.....88
 - Subscription tiers.....88
 - Burst purchases - the escape hatch88
 - Why we price the way we do90
 - Unit economics.....90
 - Path to platform economics91
- Defensibility and Moat.....91
 - The substrate compounds91
 - The discipline is harder to copy than the architecture91
 - Per-user substrate is the user's own moat.....91
 - Distribution and the trust premium92
- Traction and Validation92
 - Technical validation.....92
 - Architectural validation92
 - What we are validating next.....92
- Platform93
- Infrastructure and Security93
 - Security94
 - Roadmap.....95
 - Capital strategy95
- Innovations96
- What we Innovated97
- Lessons Learned99
- What we've learned99
- Honest Limits, Glossary, Contact104
- What We Do Not Know Yet.....105
- Glossary.....106
 - Cognition substrate106
 - KRL (Knowledge Representation Layer)106
 - Signal.....106
 - Event106
 - Meta-Context106
 - Context106
 - Pattern106

Cognition orchestrator.....	106
Cognition policy.....	106
Composer	106
Provenance.....	107
Trace	107
Synthetic cognition.....	107
Arc detector	107
Selective intelligence	107
Be Understood.....	107
Contact and Closing.....	108

A note before you read this

This started as a side project. It still feels like one in the best way.

I am writing this in the first person because I want you to know who is behind MetaOpAI. My name is Tony S. I am a security engineer with over 15 years of experience in cloud and infrastructure security, and 3 years of AI security governance at a global financial firm.

Like many people, I started using ChatGPT as a personal journal. I was exploring my own patterns, trying to improve areas of my life that needed improvement. It was incredibly helpful, until I kept hitting the same frustrating wall: the session would grow too large, the context window would collapse, and I would lose everything. All the continuity, all the hard-earned insights, gone.

I noticed several problems at scale: token bloat, inefficient full-session replay, increasing signal-to-noise degradation as sessions grow, and what I call regenerative context inflation, where the model continuously summarizes, rewrites, and appends its own generated interpretations back into future context. Over time, important information from the user becomes diluted, compressed, or lost entirely. The result is familiar to most users: once the context window collapses under its own weight, continuity breaks, sessions reset, and you're forced to redeploy important information repeatedly.

I tried other models. Same problem. Around December 2025, I got increasingly frustrated with constantly copying and pasting old context just to continue a thought. I checked the App Store thinking someone must have already solved this. Nothing came close to what I actually needed.

So I did what any sufficiently annoyed person would do: I decided to build it myself.

MetaOpAI started as a vibe-coded side project in late February 2026. I drew a rough wiring diagram on paper and began experimenting. At first I thought I'd just build a better wrapper, but

I quickly realized I had simply recreated the exact same problem I was trying to escape. That was a humbling moment.

What followed was months of obsessive work: weekends, days off, late nights, and even time on my commute. I had no deep background in AI engineering when I started, just frustration and a clear problem I wanted solved. I spent a surprising portion of my life trying to build something that did not exist yet: real, persistent continuity for personal reflection.

This unified whitepaper is the result of that journey. It's written by a guy who just wanted his AI journal not to forget everything after a few weeks, and ended up building the cognition substrate documented within the Unified Whitepaper.

The best way to describe MetaOpAI: I'm that guy in the garage trying to turn a beat-up Honda Civic into something that could gap a Koenigsegg. I was learning the AI stack as I built, and that outsider position became useful: I was not constrained by the assumptions that most AI products inherit.

I'm not claiming to have built something impossible. I'm simply saying I explored a direction most in AI don't prioritize, and what I discovered felt worth documenting.

If any of this resonates with you, the contact at the end is real. I'd rather have a small number of thoughtful conversations than a wide spray of marketing.

— Tony S.

How to Read This Paper

This paper combines what were previously two separate documents: a technical whitepaper and an investor whitepaper. We merged them because the audiences increasingly overlap, investors want to understand the architecture, and engineers want to understand the business case. More importantly, the system itself is one coherent idea, and it deserved one coherent document.

The paper is structured to support different reading paths.

If you are an investor

Read the Executive Summary, *what is MetaOpAI?* and Part I first. Then jump to Part IX (Vision and Business). The architecture sections are there if you want deeper technical detail, but the core thesis stands on its own.

If you are an engineer, researcher, or product person

Read it in order. Parts I, II, III, IV explain the design philosophy and architecture. Part IX covers innovations. Part X covers the operational lessons, including the bugs and failures that shaped the system.

If you are a prospective collaborator or curious user

Read the, *what is MetaOpAI?* Executive Summary, and Part I and II. Then jump to Part IX (Vision and Business). The architecture sections are there if you want deeper technical detail, but the core thesis stands on its own. This paper is written humbly because the work demanded it. We have been wrong in useful ways, shipped real systems, observed them under real usage, and allowed those observations to reshape our assumptions. The work is still evolving, and we expect that process to continue.

What is MetaOpAI

THESIS MetaOpAI combines a private journal with a signal intelligence engine. You log what's happening in your life, and the system extracts structured signals, behaviors, events, emotions, and relationship dynamics, then correlates them over time to reveal meaningful patterns.

MetaOpAI is a private journal combined with a signal intelligence engine. You narrate what's happening in your life, observations about yourself, your relationships, or your environments. Instead of treating these entries as disposable chat messages, the system extracts structured signals: behaviors, events, emotional framing, and relationship dynamics. These records are persisted into a governed cognition substrate we call the KRL (Knowledge Representation Layer). The KRL is organized by clear domain scopes of the user narrations:

- USER – yourself
- ENTITY – specific people
- SPACE – environments (work, family, friend groups, etc.)
- RELATIONSHIP_PAIR – dynamics between you and another person

This creates persistent continuity, turning scattered moments into a coherent, compounding record over months and years. As the KRL grows, MetaOpAI recognizes meaningful patterns

that are hard to see in the moment, trust rebuilding, withdrawal cycles, repair-and-relapse loops, shifts in reciprocity or communication, and surfaces them with evidence, confidence, and provenance rather than judgments. MetaOpAI is the private, structured, continuous space designed to help you make sense of your actual life. It sits in the gap between:

- A partner (context but no neutrality)
- A friend (neutrality but no confidentiality)
- A therapist (expertise but limited scale)
- A general-purpose AI (intelligence but no real memory)

At its core, the real innovation is the KRL itself as this assists us with the persistency. The language model is treated as a rented, swappable tool at the edge. The durable value, memory, patterns, and understanding, lives in the governed KRL that becomes more valuable the longer you use it. The company motto is Be Understood.

Executive Summary

THESIS Large language models are becoming commodities. The real moat in consumer AI is the cognition substrate, the layer between the model and the experience that actually remembers, recognizes patterns, and connects insights over time. MetaOpAI is building that substrate, focused on the three areas where humans carry continuity: yourself, the people in your life, and the environments you inhabit. The product is signal intelligence, our motto is Be Understood, and the business is a subscription that monetizes the deep, compounding experience of being truly seen over time.

MetaOpAI is a cognition engine optimized for signal intelligence. We focus on the three areas where humans naturally carry continuity: yourself, the people in your life, and the environments you inhabit. We are not a chatbot, a therapy product, or a generic journaling tool.

MetaOpAI is a private journal paired with a structured cognition engine. You narrate what is happening in your life, and the system extracts structured signals, behaviors, events, emotional framing, and relationship dynamics. These signals are stored in the Knowledge Representation Layer (KRL), organized by clear scopes:

- USER
- ENTITY
- SPACE
- RELATIONSHIP_PAIR

As the KRL grows, MetaOpAI surfaces meaningful patterns over time, trust rebuilding, withdrawal cycles, repair-and-relapse loops, shifts in reciprocity, changes in communication,

and emotional state transitions, all with evidence, confidence, and provenance rather than judgment.

The system functions less like chat memory and more like long-term behavioral observability, enabling true continuity across months and years.

The language model itself is treated as a rented, swappable tool at the edge. The real value lives in the KRL, the persistent, versioned cognition substrate underneath.

MetaOpAI is built around five product-facing cognition layers:

1. Journal
2. KRL (Knowledge Representation Layer)
3. Pattern Engine
4. Cognition Orchestrator
5. Composer

Together, these layers create an AI experience that becomes meaningfully more valuable the longer it is used.

Our motto is:

“Be Understood.”

The deep experience of being truly seen and understood over time.

PART I

Foundations

What we are trying to do, why we think it matters, and the design philosophy that has guided every architectural decision.

A Familiar Moment

You've been here. Trying to make sense of someone.

A friend who is pulling away. A manager whose tone shifted last month. A partner whose evenings feel a little different lately. You sense something, but you cannot quite name it. So, you reach out for help.

MetaOpAI is built around a simple observation: people often sense that something is changing long before they can clearly explain it. A friend slowly pulls away. A manager's tone subtly shifts over the course of a month. A partner's routines begin to feel emotionally different. Individually, these moments seem small or ambiguous, but together they form patterns. Human life produces a constant stream of emotional and behavioral signals, changes in communication, reciprocity, consistency, tension, enthusiasm, avoidance, presence, and absence, yet most tools treat these moments as isolated events instead of connected data points unfolding over time.

When people try to make sense of these situations, they usually turn to conversation. They talk to a spouse, a friend, a therapist, or increasingly an AI chatbot. But each option breaks down in a different way. Partners have context but are emotionally involved. Friends may be neutral but introduce bias, exposure, or incomplete understanding. Therapists provide expertise but only see small snapshots separated by days or weeks. General AI chatbots may be intelligent, but they lack continuity: each conversation eventually collapses into a fading context window, making durable pattern recognition impossible.

MetaOpAI was designed specifically for this gap. Instead of treating conversations as disposable chat sessions, the system treats them as long-term signal intelligence. The platform functions like a private journal connected to a cognition substrate that continuously extracts structure from narration: behaviors, events, emotional framing, relationship dynamics, and environmental shifts. Those records are persisted over time across four dimensions, yourself, other people, environments like work or family, and the relationships between individuals. The result is not simply memory, but continuity: a system capable of correlating signals across weeks, months, or years rather than within a single session.

Over time, MetaOpAI begins surfacing patterns that are difficult to see from inside your own life while you are emotionally embedded within it. Trust rebuilding slowly. Withdrawal accelerating. Repair-and-relapse cycles repeating. Changes in reciprocity, tone, reliability, or emotional state emerging gradually over time. Importantly, the system does not present these as verdicts; it presents them as evidence-backed observations with confidence and provenance. Underneath the interface, the core innovation is the cognition substrate itself, a structured continuity layer positioned between the user and the language model. The language model is treated as a

replaceable commodity at the edge, while the actual value compounds inside the persistent substrate: the memory, the patterns, the relationships, and the accumulated understanding of a person’s life over time.

Today, you have four choices. None of them is built for this.

A FAMILIAR MOMENT

You've been here.

Trying to make sense of someone.

A friend who's pulling away. A manager whose tone shifted. A partner whose evenings feel a little different. You sense something — but you can't quite name it. So you reach out for help.

Today you have four choices. None of them is built for this.

OPTION 01 — A PARTNER OR SPOUSE

**They love you.
They're missing the backstory.**

— THE GAP · CONTEXT TAX & WELL-MEANING NOISE

You narrate from the beginning — personalities, history, baseline, deviation. By the time you reach the question, half an hour is gone. The feedback is loving and well-meaning — but downstream of a story you had to rebuild.

OPTION 02 — A FRIEND

**They have an opinion.
Or you can't tell them at all.**

— THE GAP · BIAS, EXPOSURE, CONFIDENTIALITY

Friends have their own emotional stake. Opinions form. Sides get taken. Some situations are too sensitive, too professional, or too close to home to confide in anyone in your circle in the first place.

OPTION 03 — A THERAPIST

Once a week. About you.

— THE GAP · COST & COVERAGE

\$150–\$300 a session. Once a week if you're lucky. Therapy is built to help you understand your internal experience — not to keep structured memory of every person in your life and the patterns between them.

OPTION 04 — AN AI CHATBOT

Smart. Willing. Amnesiac.

— THE GAP · NO MEMORY, NO PATTERNS

You open ChatGPT or Claude, narrate the situation, get something insightful. Tomorrow, blank slate. The next event gets analyzed without yesterday's context. Patterns cannot compound if memory cannot persist.

MetaOpAI is the fifth option.

Private. Structured. Continuous. Built for the in-between hours of your actual life.

Figure A, *The four-option problem. A partner has context but no neutrality. A friend has neutrality but no confidentiality. A therapist has expertise but no scale. An AI chatbot has scale but no memory. MetaOpAI is built to be the fifth option.*

Option [1]: A partner or spouse

They love you. They are missing the backstory. You start narrating from the beginning, the personalities, the history, the baseline, the deviation, finally the event itself. By the time you reach the question, half an hour is gone. Their feedback is loving and well-meaning, but it is downstream of a story you had to rebuild from scratch. You walk away supported. Not clearer.

Option [2]: A friend

They have an opinion. Or you cannot tell them at all. A friend brings their own emotional stake. Opinions form, sides get taken, things get repeated. And some situations are simply too sensitive, too professional, or too close to home to confide in anyone in your circle in the first place.

Option [3]: A therapist

Once a week. About you. \$150 to \$300 a session, once a week if you are lucky. Therapy is built to help you understand your own internal experience, not to keep structured memory of every person in your life and the patterns between them. It is a precious tool. It is not this tool.

Option [4]: An AI chatbot

AI today is smart and helpful, but forgetful. You explain your situation, it gives insight, and over time the conversation grows into repetitive context and noise. Eventually the session becomes too large, breaks down, or you start a new chat and lose continuity. The next event gets analyzed without the history that gave it meaning. Patterns cannot compound if memory does not persist.

MetaOpAI is the fifth option. Private. Structured. Continuous. Built for the in-between hours of the actual life you are trying to make sense of.

The Anomaly We Noticed with AI Wrappers and Gen Purpose AI (Chat)

Every general-purpose AI chat, ChatGPT, Claude, Grok, DeepSeek, Gemini, works the same way under the hood. The way it works has costs that compound.

The pattern observed is simple: the user types X, the AI responds with Y. The user then types D, and the AI answers using X, Y, and D because most AI systems and wrappers re-send prior conversation context back to the model each turn. The user types F next, and now the model sees X, Y, D, B, and F. The context window keeps growing. Tokens compound. Noise increases. Over time, the model can end up processing more AI-generated synthesis than original user input.

It feels like the AI gets overwhelmed and starts drowning in its own memory.

This is not a bug. Large language models are fundamentally stateless. They do not naturally remember previous conversations between calls, so providers re-send prior context to create the illusion of memory. That explanation is broadly correct, with one important nuance: modern systems often use additional techniques like summarization, retrieval, caching, memory stores, and truncation instead of literally sending the entire conversation forever. But the core observation is valid: continuity in most AI systems is simulated through context reconstruction, and that creates scaling, cost, and signal-quality problems over long sessions.

From a business perspective, many AI wrappers face a difficult economics problem because API cost is tied to token usage. If the application keeps re-sending growing conversation history, even compressed, operational costs increase significantly over time while context quality can degrade as the window fills with repetitive or synthesized information.

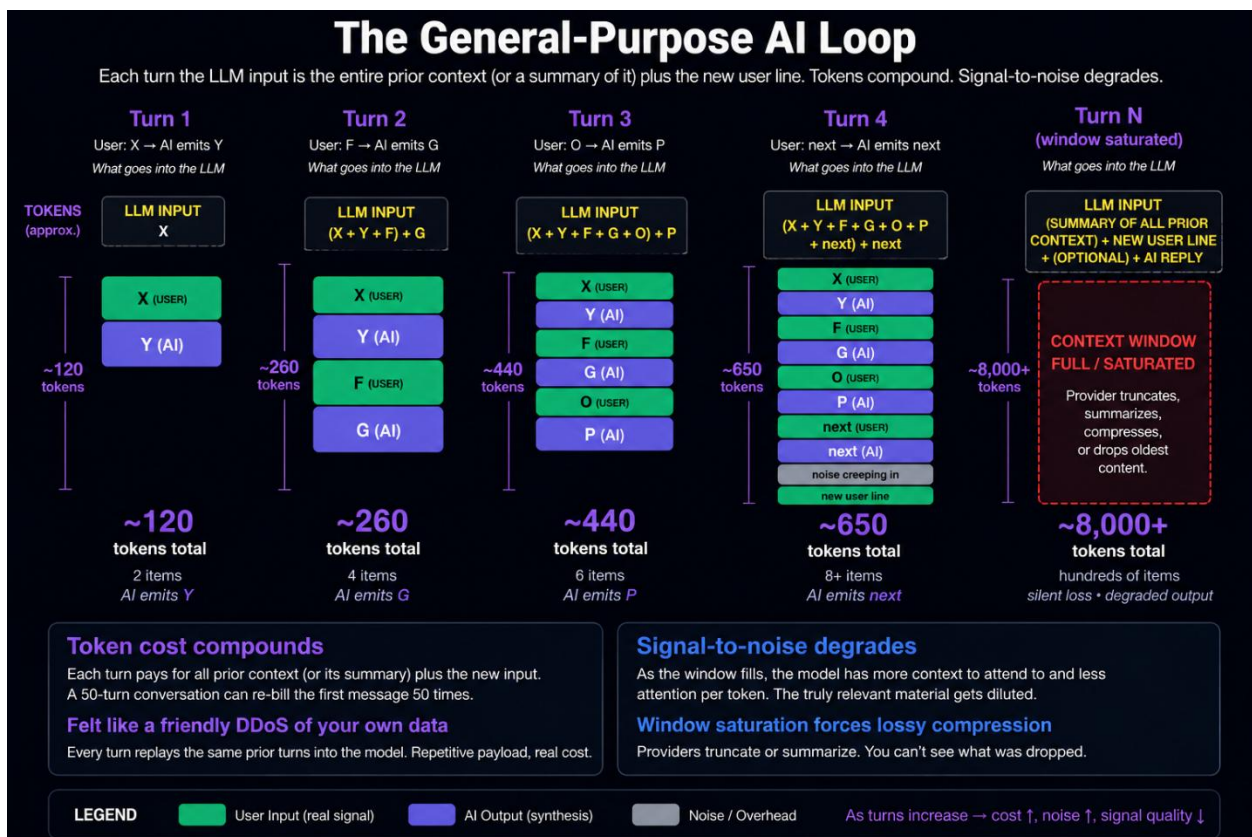


Figure B, The general-purpose AI loop. Every turn replays the full history into the model. Token cost compounds. Eventually the window saturates and providers silently truncate or compress, losing fidelity.

The costs nobody talks about

Replay-heavy AI systems carry costs that compound over time. Most AI wrappers inherit how chat LLMs work — replaying a growing conversation history into the model every turn. As the context window expands, token costs climb, and eventually it stops being practical.

To be precise: against a naive system that replays the full transcript with no caching, the gap is dramatic. Against a well-cached wrapper — the realistic competitor — it's smaller, a few times rather than tens.

But the multiplier was never the point. The point is the shape of the curve. A replay system's per-turn cost rises with conversation length; the only question is how steeply. MetaOpAI's per-turn cost stays flat no matter how long the relationship runs.

Caching can flatten a competitor's curve. Only architecture makes it constant.

MetaOpAI also has meaningful token usage, but because its architecture retrieves a small, scope-correct slice rather than replaying history, per-turn cost is effectively constant — $O(1)$ in conversation length — instead of growing with it. That constancy, not any specific multiplier, is the durable economic claim. With further optimization from experienced AI and systems engineers, those costs could theoretically be reduced even further.

- **Cost to the app owner:** Re-sending large conversation history every turn increases API costs because more tokens are processed. Most apps pass those costs to the user through subscriptions, limits, or pricing. For LLM providers, bigger context windows often mean more revenue.
- **Cost to the user:** As the conversation grows, the context window fills with old summaries and AI-generated synthesis. The noise-to-signal ratio increases, and the AI can feel overwhelmed or less focused. Users often respond by typing even more detail to correct the AI, which adds even more noise. The more material the model has to attend to per token, and the less attention it can give to the parts that actually matter. Eventually the window fills entirely. At that point, providers truncate or summarize the oldest turns. The model cannot tell you what was lost, because from its perspective, those tokens simply do not exist anymore. You feel it as the AI forgetting things, drifting off-topic, or repeating itself. The technical name is context window saturation. The felt experience is talking to a friend whose attention is being slowly diluted by their own memory.
- **Cost to energy and infrastructure:** Processing massive conversation histories requires more GPU compute, memory, and power consumption. Even with compression and caching, replaying long sessions creates significant infrastructure overhead because the system is repeatedly rebuilding continuity instead of truly remembering it.

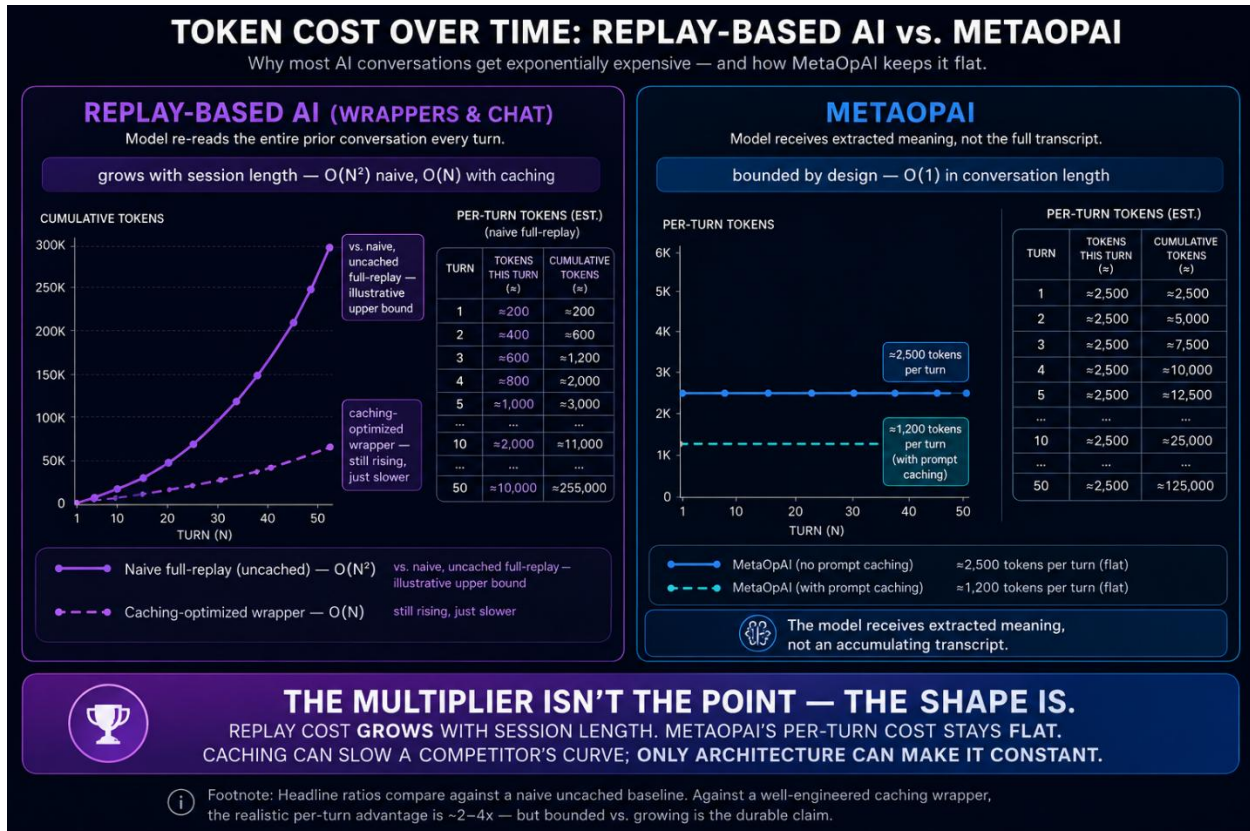


Figure C — Why replay-based context cost grows while MetaOpAI's stays flat. In a replay architecture, every turn re-sends prior context — both the user's input and the model's own prior output — back into the model, so per-turn input tokens rise turn over turn. The headline ratios compare against a naive, uncached full-replay baseline and represent an illustrative upper bound; a well-engineered wrapper using prompt caching grows more slowly, but still grows. MetaOpAI's per-turn cost is bounded by design — roughly 2,500 tokens per turn (≈1,200 with caching) regardless of session length — because the model receives extracted meaning, not an accumulating transcript. The durable claim is not a multiplier but a shape: replay cost rises with conversation length; ours stays constant.

Why the wrapper consensus does not fix it

The wrapper model is profitable for the LLM provider, but largely indifferent to the user, the application owner, and the investors funding those applications. To solve the bloat problem, you have to stop treating the chat session itself as persistent memory and build a different foundational layer underneath. That is what MetaOpAI aims to be, not a wrapper, but a substrate.

What We Built Instead

Instead of replaying the conversation, we extract its meaning and store it. The model receives meaning, not transcript.

The idea is simple enough to write in one sentence: every time a user narrates something, we run extractors that pull structured records out of it, signals, events, context, meta-context, and store them in a substrate organized by who or what the record is about. Future turns retrieve only the small, relevant subset for the current scope. The model gets a compact, dense prompt. The window stays small. The bill stays sane. The continuity grows.

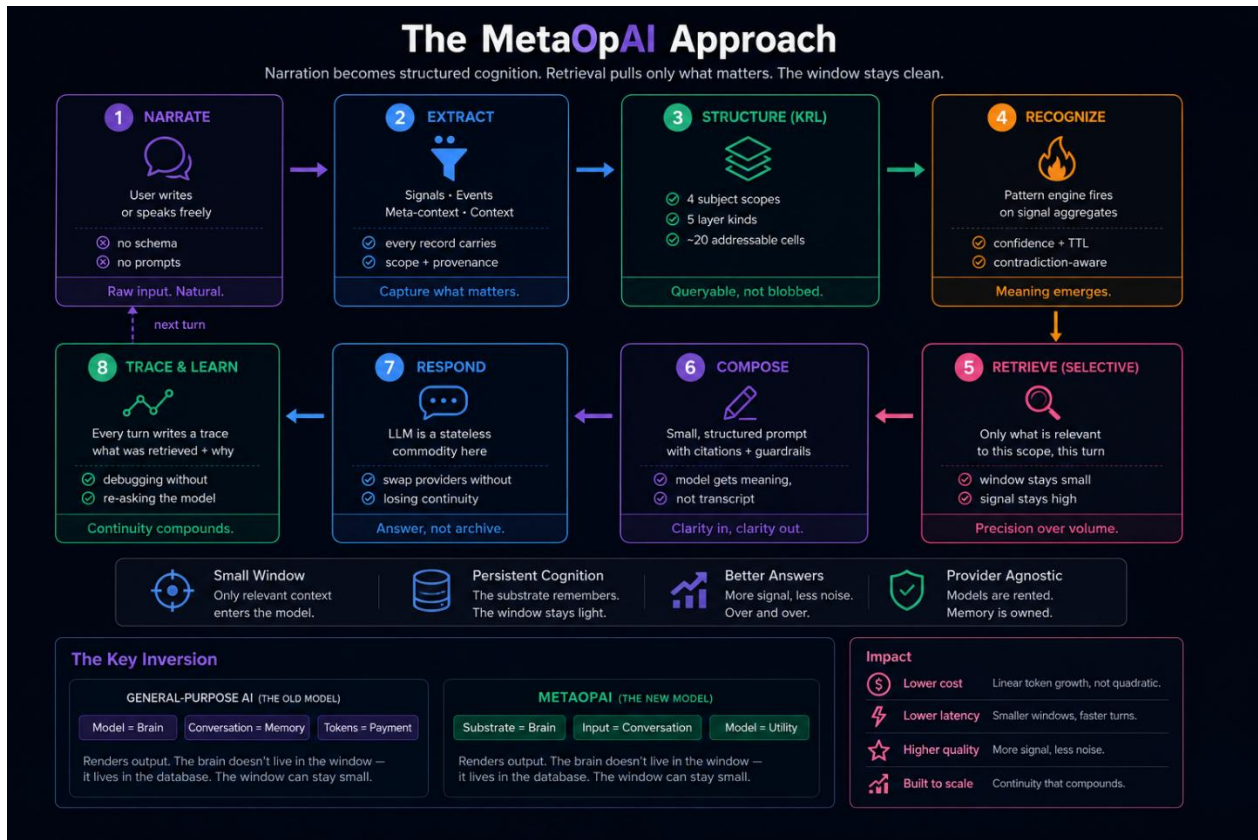


Figure D, The MetaOpAI flow. Narration is extracted into structured cognition. Retrieval pulls only what is relevant to the current scope. The LLM becomes a stateless commodity at the rendering edge.

We organize the substrate as a small cognitive matrix: four subject scopes, user, entity, space, and relationship pair, across five-layer kinds: signals, events, meta-context, context, and patterns. That creates twenty addressable cells instead of one giant conversational blob. Every record carries provenance: which extractor produced it, from which message, at what time, and at what confidence.

The subject scopes, users, entities, spaces, and relationship pairs, are the dimensions of the system, while the layers, signals, events, context, meta-context, and patterns are the properties that exist across those dimensions.

Above that sits the pattern engine. It does not react to isolated moments; it waits for evidence to accumulate across time and scope before recognizing recurring dynamics like avoidance loops, repair-then-relapse cycles, reciprocity drift, or trust trajectories.

We did not invent the idea of memory in AI. The broader industry is already moving toward long-term memory systems. Where MetaOpAI differs is in the discipline around what gets stored, what gets surfaced, and what gets suppressed. The same instinct that led us to challenge conversational token bloat shaped the substrate itself: if you store everything, you simply move the bloat from the context window into the database. The goal is not accumulation. The goal is selective continuity. The conversation is not the memory. The substrate is.

Why we did not build a wrapper

We could have. A wrapper is faster, cheaper, and easier to demo. You take an LLM, drop a memory feature on it via the provider's APIs, and ship. We chose not to. The wrapper path locks you into one model provider, one memory implementation, and one pricing curve. More importantly, it inherits the same architectural assumption, that the session is the memory, and bolts a workaround on top of it. We wanted to question the assumption.

So, we built the substrate ourselves. Postgres for durable cognition, JSONB for the structured records, a pattern engine that fires on aggregates, a cognitive orchestrator that decides what to retrieve, and a trace layer that lets us see what the system did and why. Every layer is independently versionable. We can swap the model provider, change the extractor prompts, retune the retrieval policy, without changing the others. The architecture is the moat. The model is the rented commodity.

PART II

Product Experience

How the ontology shows up in the UI, onboarding, journals, analytics, the surfaces a user actually touches.

Product Tour

Onboarding. The screens. The journals. The analytics. The settings. What every surface does and why it exists.

The core investigation path is simple: explore patterns, map where they live, inspect how they formed, and verify the evidence.

This chapter walks through MetaOpAI as a user would meet it. The aesthetic is locked, the design came from a partner who did the work properly, and the architecture lives behind every screen. We document the surfaces here so a reader can see the product as a coherent whole, not a list of features.

Onboarding

Onboarding is deliberately short. The identity form captures the small set of stable facts that the substrate uses to ground future cognition, primary role, cultural background, communication style, and routes the user into creating their first space. There is no AI chat during the identity step; the only AI interaction in onboarding is a capped two-turn context conversation that runs after the first space is created, where the system asks what the user wants to track and listens. The synthetic cognition layer writes a small low-authority blob from those turns to help the substrate bootstrap. As real continuity accumulates, the synthetic blob's weight tapers.

After onboarding, the tutorial offers a guided pass over the major surfaces. It is opt-in and skippable. We did not want a forced tour blocking the product.

The investigation workflow has four layers: KRL Explorer answers what patterns exist; Graph Explorer answers where those patterns live across people and spaces; KRL Formation answers how a pattern was constructed; Evidence Chain answers what exact observations support it.

The KRL Explorer

KRL Explorer. The KRL Explorer is the cockpit landing pane and the entry point into your Knowledge Representation Layer. It surfaces your visible pattern bank as a scannable, queryable index — every pattern the cognition substrate has detected across your spaces, entities, and self-reflection journals shows up here as a card carrying its humanized label, lifecycle state (active / changing / repeating / weakening / needs-review), confidence score, evidence count, scope (which subjects it applies to), and last-seen timestamp. The top of the pane gives you two complementary investigation modes: a Structured Query terminal where you filter by entity, space, lifecycle, confidence floor, pattern type, cross-space-only, needs-review-only, and time window with CSV export — and an Ask AI box that takes a natural-language question and routes through the analytic voice to answer in conversational form over the same substrate. The page is the "what is happening in my life right now, by the numbers" view — emotionally light, analytically dense, designed to be opened the way someone opens a Bloomberg terminal at the start of a workday.

KRL Explorer
Explore Understand Evolve

Search patterns, people, contexts, signals, events...

Structured query Ask AI

SEARCH MODES

- Natural Language (Ask anything)
- Structured Query (Build a filter query)
- Graph Explorer (Explore connections)

QUICK SEARCHES

- My top patterns
- Recently strengthened
- Most observed
- Cross-space patterns
- Contradicted patterns
- Weak / fragile patterns

RECENT SEARCHES

No recent searches yet.

Search Tip

Use natural language to ask complex questions about your cognition.
"Why do I shut down around Mike?"

Explore Patterns
Find and explore cognition patterns across your knowledge graph.

Domain: All Scope: All Sort by: Confidence Show: 10

PATTERN	DOMAIN	FORM	CONFIDENCE	OBSERVED	TREND	LAST SEEN
Unreliability pattern Work	Trust	Cluster	1.00	7x	⚠ Weakening	May 23, 2026
Dismissiveness pattern School + Friends	Communication	Cluster	1.00	4x	➔ Strengthening	May 16, 2026
Repair-then-relapse cycle Relationship	Conflict	Cycle	0.92	8x	➔ Strengthening	May 24, 2026
Trust & reliability pattern Work	Trust	Composite	0.80	8x	⚠ Weakening	May 23, 2026
Unreliability pattern You	Trust	Cluster	0.80	5x	⚠ Weakening	May 16, 2026
Contextual Divergence Pattern Work + Friends	Pattern	—	0.75	2x	➔ Strengthening	May 28, 2026
Anchor Space Pattern	Pattern	—	0.75	2x	➔ Strengthening	May 28, 2026

ACTIVE FILTERS Clear all
None

FILTER BY

ENTITIES
e.g. Mike

SPACE
e.g. Work

MIN CONFIDENCE
≥0.5 ≥0.6 ≥0.7 ≥0.8

Cross-space only

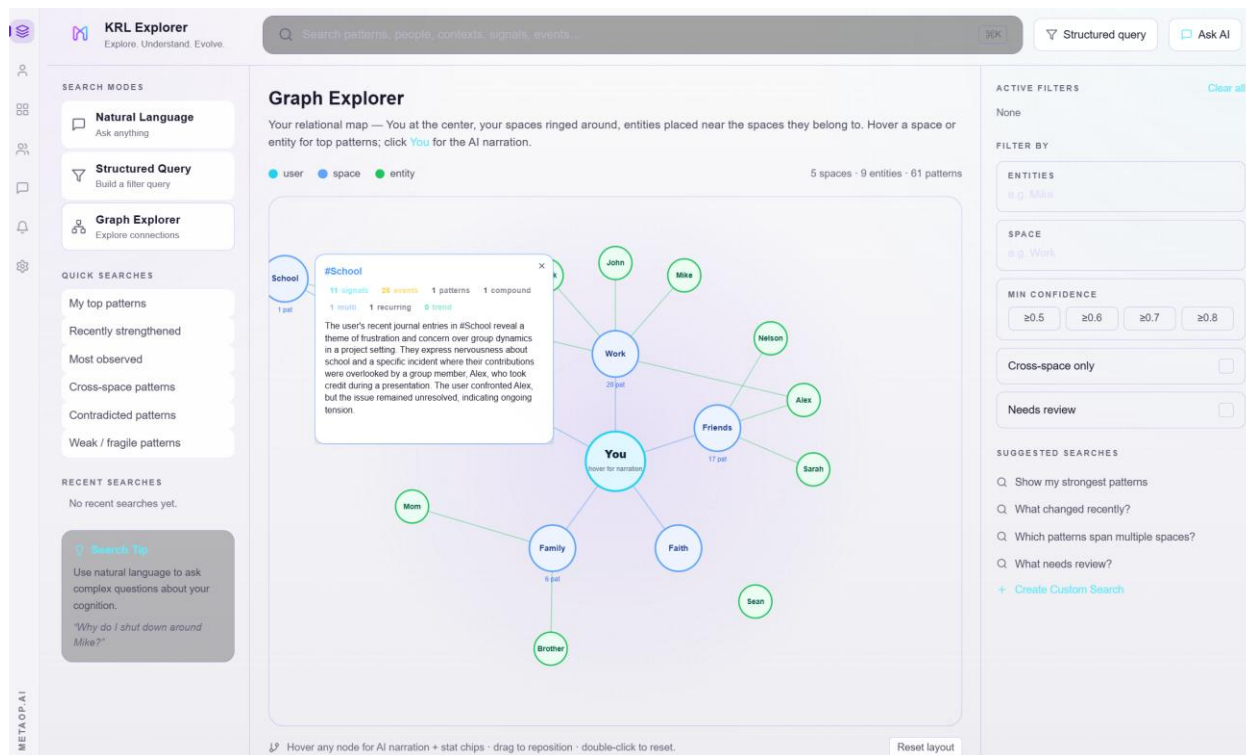
Needs review

SUGGESTED SEARCHES

- Q Show my strongest patterns
- Q What changed recently?
- Q Which patterns span multiple spaces?
- Q What needs review?
- + Create Custom Search

The Graph Explorer

Graph Explorer. The Graph Explorer is the spatial counterpart to the Explorer's list view — it shows the same substrate as a relationship topology rather than a feed. The layout is concentric: You sit at the center, your spaces orbit in the inner ring, and your entities sit in the outer ring connected by their space memberships, so a single glance tells you who is in your life, where they live (which spaces), and how the spaces relate to each other. Hovering or clicking any node opens a popup with seven deterministic KRL counters — signal count, event count, pattern count, compound-pattern count, multipattern count, recurring-pattern count, and trend count — pulled directly from the `node_counts.py` aggregation path against the substrate (no name-match heuristics), plus an AI-narrated paragraph describing what that person, space, or you-as-subject actually look like across the data. The graph is interactive: you can drag any bubble to reposition it, click to pin a focus node and dim the rest, and the popup persists so you can read the narration. This is the "where is everything connected, and what does each node carry" view — the investigative complement to the Explorer's list.



The KRL Formation

KRL Formation. The KRL Formation page is the "why does this pattern exist?" view for any single pattern you click into from the Explorer or Graph. It opens the pattern's construction record and lays it out as a 4×4 matrix: four dimensions down (USER / ENTITY / SPACE / RELATIONSHIP_PAIR) by four properties across (Signal / Event / Context / Meta-Context), with each populated cell showing the contributing extracted record, the quoted source text, the confidence the extractor assigned, and how many times that observation repeated. This is the architecture of how the cognition substrate composed the pattern — not the conclusion, but the raw multi-subject, multi-property compounding that produced it. Above the matrix sits the Genesis Sequence: a horizontal strip of circles showing the chronological build-up of the pattern from first signal through current state, so you can see whether the pattern formed quickly from a dense cluster or accumulated slowly across weeks. The page exists because the user is paying for the pattern engine's reasoning, not just its output — Formation makes the reasoning legible.

Dismissiveness pattern · Why does this exist? · KRL Formation

How "Dismissiveness pattern" was formed

Every pattern crystallizes from the Knowledge Representation Layer — **dimensions** (who/what) × **properties** (kind of evidence). Lit cells fed this pattern, dashed cells held nothing.

Domain: Communication | Form: Cluster | Scope: Entity

Tracked 112+ across 9 KRL elements — still recurring

WHAT THIS PATTERN MEANS
Repeated instances of being brushed off, minimised, or talked over, observed close together in time.

GENESIS SEQUENCE

- Trigger**
Alex made fun of Sarah's engagement ring, which may have been perceived as rude and dismissive of her feelings.
- Noticed**
Alex made fun of Sarah's engagement ring, which may have been perceived as rude and dismissive of her feelings.
- Interpreted**
Detected: Dismissiveness pattern on Alex — 4 contributing signals over a 14-day window. No "At"
- Recurrent**
4 observations reinforced it.
- Spread**
2 spaces: School, Friends.
- Confirmed**
— 1.00
Calibrated from the evidence.
- Crystallized**
active
Dismissiveness pattern

	SIGNAL	EVENT	CONTEXT	META-CONTEXT	-- PATTERN
User You	SIGNAL User expresses feeling tired today, indicating a low energy state. User expresses feeling tired today, indicating a low energy state. tracked 79+	EVENT User is officially the new Director of Security Engineering starting Monday. Just finalized — I'm officially the new Director of Security Engineering starting Monday. tracked 8+	CONTEXT Persian-American I'm Persian-American tracked 7+		<p>Dismissiveness pattern</p> <p>COMMUNICATION</p> <p>Confidence 1.00 Strengthening</p> <p>Show full evidence table</p>
Entity Entity	SIGNAL Alex took all the accolades for the group project and didn't accredit the user for the work they did. Alex took all the accolades for the group project and didn't accredit the user for the work they did. tracked 9+		CONTEXT first_only @Alex tracked 1+		
Space School					
Space Friends	SIGNAL User indicates that family dynamics contribute to stress during significant events. User indicates that family dynamics contribute to stress during significant events.	EVENT Sarah and Alex joined user for poker night. @Sarah and @Alex joined for poker night	CONTEXT close		

The KRL Chain of Evidence

Evidence Chain. The Evidence Chain is the deepest layer of the investigation flow and the audit trail behind every pattern claim the system surfaces. Where Formation shows the matrix of contributing record types, the Evidence Chain traces the pattern all the way back to the specific source records that fed into it — the exact journal sentences you wrote, the events the system extracted, the meta-context fragments it persisted, each with its date, its original quote, the conversation it came from, and the dimension and property slot it landed in. The chain is presented chronologically so you can follow the cognition substrate's reasoning step by step: this signal on this date, then this event two days later, then this contradicting signal, then this reactivation. This matters because MetaOpAI is positioned as signal intelligence rather than therapy, and signal-intelligence output requires audit-grade provenance — a user has to be able to ask "show me the receipts" and get the actual underlying observations. The Evidence Chain is the receipts.

The screenshot displays a web interface for an evidence chain. At the top, there is a breadcrumb trail: "Dismissiveness pattern" > "KRL Formation" > "Evidence Chain". To the right of the breadcrumb are "Back" and "X" buttons. Below the breadcrumb, the title "Full evidence chain" is followed by a subtitle: "Every observation that built **Dismissiveness pattern**, in order — each one moving confidence up or down." The main content is a vertical list of six items, each with a circular icon on the left and a date on the right. The items are:

- DISMISSIVE** (2026-05-17): Alex acted as if the situation was no big deal after the user confronted him.
- DISMISSIVE** (2026-05-17): Alex acted as if the situation was no big deal after the user confronted him.
- DISMISSIVE** (2026-05-17): Alex acted as if the situation was no big deal after the user confronted him.
- CONTRADICTION** (2026-05-10): User found the poker night cool, noting Alex almost flipped out when they won a hand bluffing, indicating a light-hearted moment. — Counter-evidence to Dismissiveness pattern: a humor_shared signal in the same window suggests the pattern is partial.
- DISMISSIVE** (2026-05-06): Alex made fun of Sarah's engagement ring, which may have been perceived as rude and dismissive of her feelings.

A vertical sidebar on the left contains several icons, and the text "METAOPAI" is visible at the bottom left of the interface.

Profile - the user surface

The Profile is the user's home. It is the read view for everything the substrate has learned about the user themselves. There is a noticeable load latency here that we own honestly, the surface is data-heavy and we have not yet finished optimizing every query path. What it shows: time invested across the product, spaces ranked by narration volume, entities ranked by volume, top five spaces and top five entities by activity. It also includes the User Intel Report, the top signals and top patterns the system has observed about the user, experimental self-reflection lenses such as communication-style and relationship-pattern estimates. We label these Experimental because they are pattern interpretations, not clinical assessments, and we want users to read them as such.

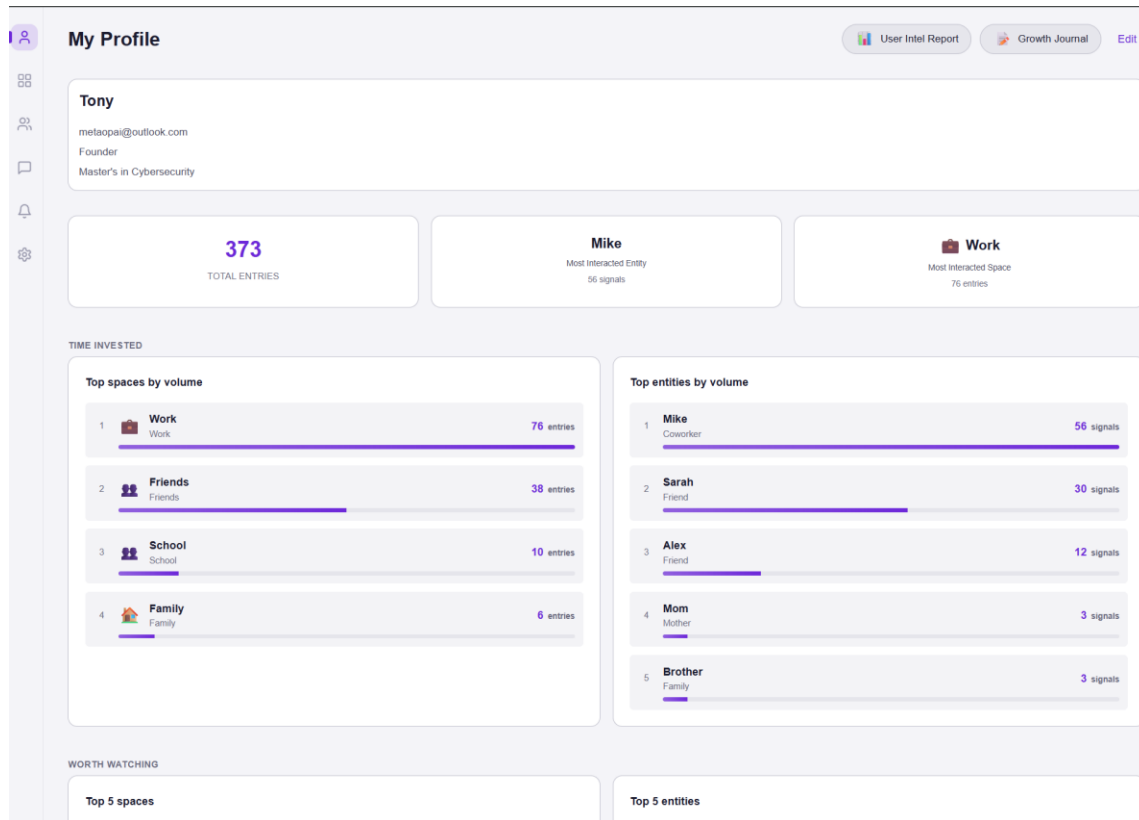


Figure E, The top of the Profile pane of our QA account using a federated auth via Microsoft. We can see journals that are most frequent and users that are most frequently mentioned. Shows total entries.

Behavioral metrics

Inside Profile we surface a behavioral panel: typing rhythm, response latency, sloppiness drift, stress state, word choices, entry length trend. These come from the behavioral capture model, they are collected only when the user has consented (the toggle is a collection gate, not a display gate), they are tied to the content they were captured alongside, and they inform the substrate's reading of context without being shown as judgments. Why we collect them: a sentence typed quickly with high error rate at 1 a.m. is different from the same sentence typed slowly at 10 a.m., and the substrate's job is to register that. We surface the inputs so users can see what we are using.

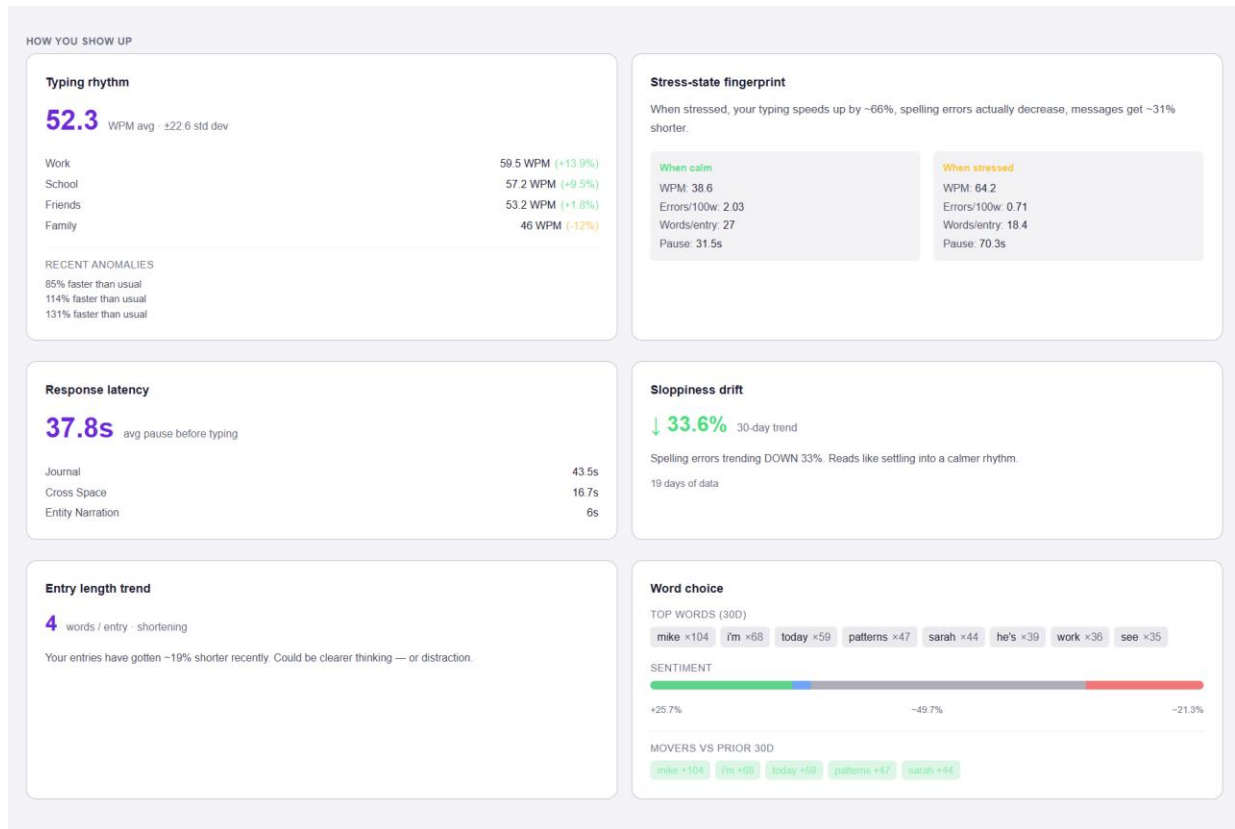


Figure F, Textual communication inherently lacks emotional nuance. To bridge this gap, this system uses optional, user-consented signals — typing rhythm, response latency, and entry texture — captured locally and only when the user has enabled them, to help the system read context more accurately. (Softened 'behavioral biometrics' / 'analyze' framing) By benchmarking these passive metrics against a user's historical norm, the platform extracts localized signal intelligence to infer cognitive load and emotional state shifts without relying solely on explicit written language this can be disabled within settings.

Activity heatmap

A simple calendar heatmap of session activity, scoped to the user. Useful for seeing rhythm and detecting weeks of avoidance or intensity. We do not assign meaning to the heatmap, it is just a visualization the user can interpret themselves.

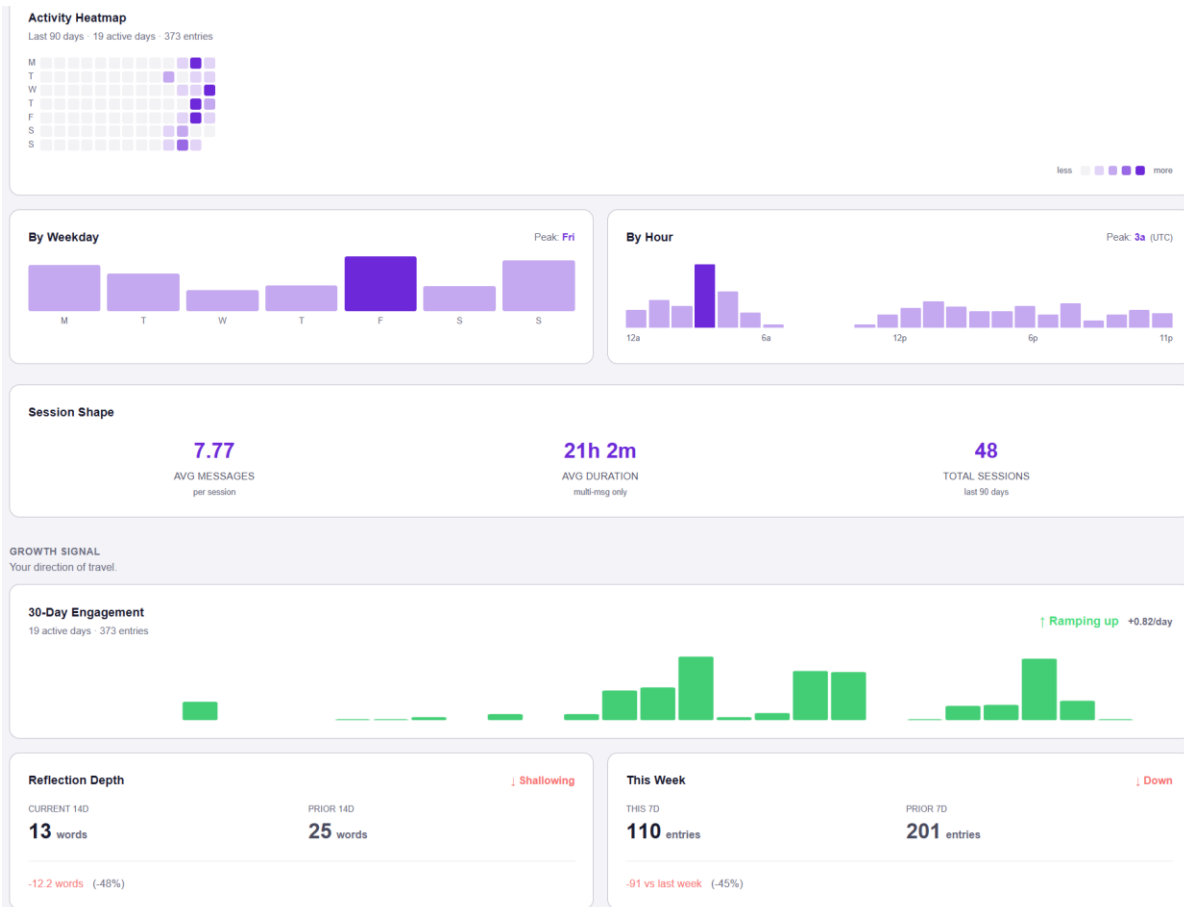


Figure G, Graphical visual of user activity and engagement within the application.

Spaces - Mirror

the systemic surface

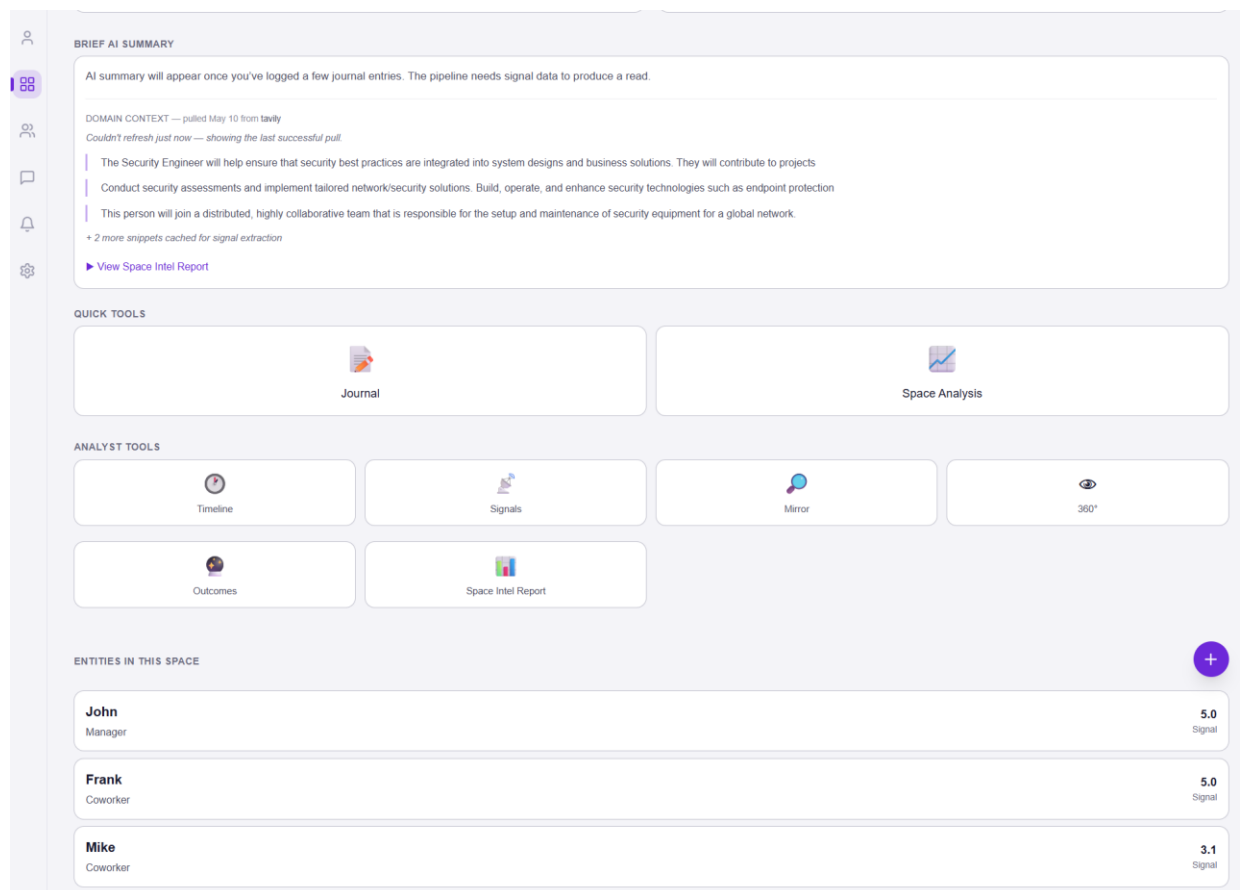
Spaces are sections of life. Work. Family. Friends. A particular relationship. A community. Each space is a journal of its environment, plus a suite of analytics that read from the same scope-correct slice of the substrate. The analogy that fits best is a composition notebook with built-in analytics: you write in it, and the notebook tells you what it has observed.

Spaces are dedicated environments inside MetaOpAI for different parts of a user's life, such as Work, Family, Marriage, or personal projects. Users assign entities (people) to those spaces, and each space develops its own continuity, context, and dynamics over time. The experience is intentionally similar to a composition notebook: each notebook represents a different area of life, while the built-in analytics and cognition tools act like the reference tools often found

in the back of a notebook, helping users better understand the patterns unfolding inside that environment.

The Space Journal

The journal is where narration happens. The user writes (or speaks), the orchestrator runs, the LLM responds with continuity-aware framing, and the substrate quietly extracts. Entities are referenceable via @mention, typing @ shows the entities assigned to this space and inserts the right reference so signals attributed to that mention land in the correct entity record. Sessions auto-archive after a quiet period or can be archived manually. Archived sessions live in the Timeline and can be resumed in place, picking up the same thread rather than starting a new conversation.



*Figure H, The Journal Home Page Architecture. A multi-journal dashboard tied to subscription levels. Selecting a journal activates a localized RAG context. The UI exposes three core functions: the **Journal Workspace** for entries and AI interaction; a **Dual Analytics Suite** for signal tracking and zero-ingress QA; and an **Entity Panel** where user profiles are invoked and updated to maintain persistent memory.*

The Space Analysis

User can interact and ask questions to the AI about all things in the space without concerning about generating signals.

Timeline - sessions over time

The Timeline shows every past session as a page-turner UI. Each page has a summary, the entities present, the signals extracted, and the option to expand the full transcript.

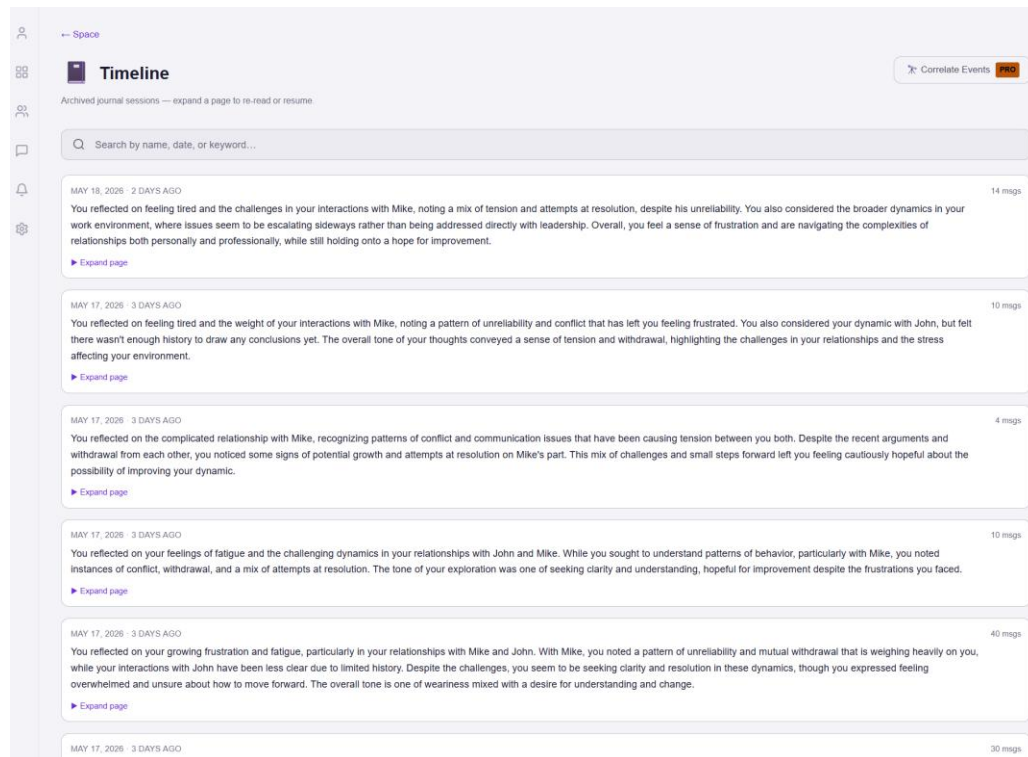


Figure 1, The Dynamic Timeline Feed. A chronological interface where sessions act as sequential journal pages. Committing a session archives an AI-generated summary; while expanding it displays the full history. Users can resume sessions to update their chronological order. Highly frequented pages are dynamically tagged by the AI as Rumination events and given a "worn" visual aesthetic, signaling deep user focus or unresolved cognitive loops on that specific entry

Timeline - Correlation Feature (Subscription Only)

Correlation features in the Timeline let users see how sessions relate, events that reference earlier events, patterns that span multiple sessions, anchors that connect back to specific moments.

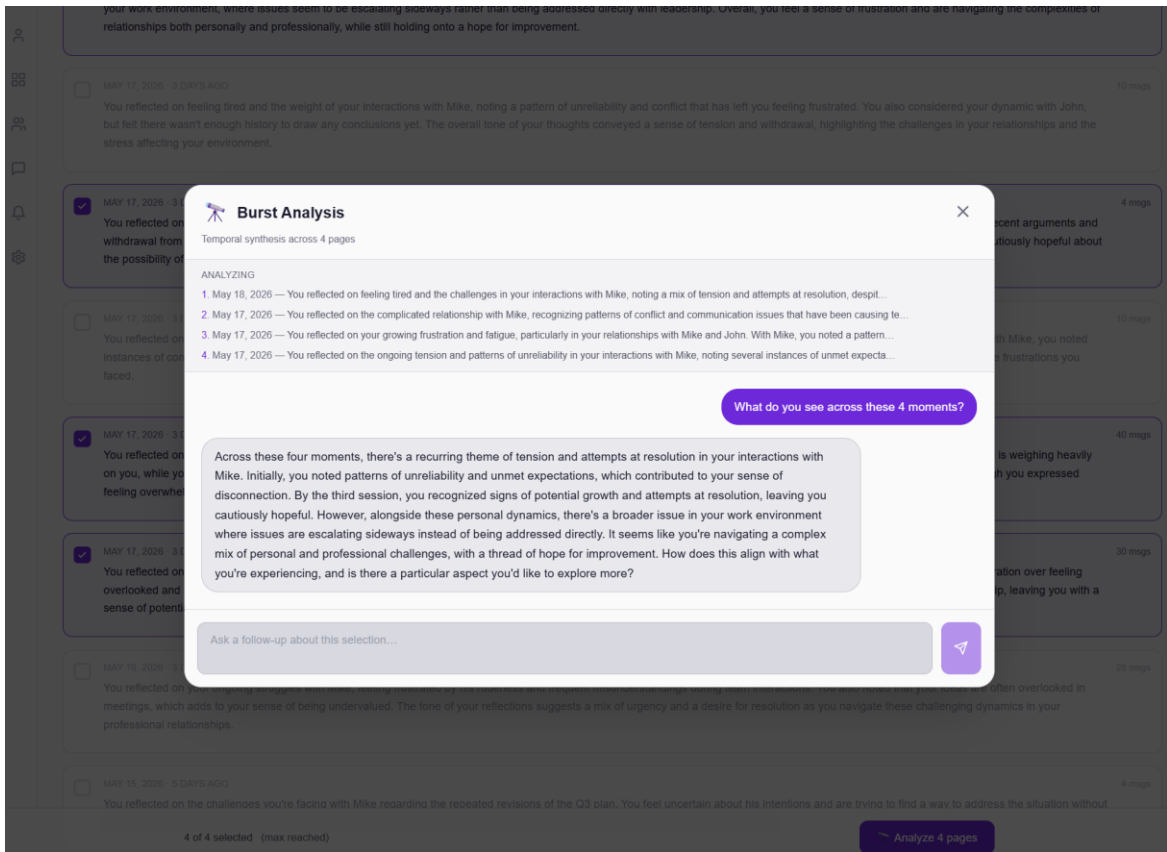


Figure J, The Correlation “Burst Analysis” the idea is correlation, where you can select multiple epochs up to 4 and it will import them into the AI window displayed above and correlate across those multiple events to look for hidden pattern, you can ask q/a and interact with the AI. This is only accessible in the subscription model.

Signals - the narration-level view (Subscription Only)

Signals show the per-utterance behavioral markers the extractors have written: withdrawal, vulnerability punished, repair_attempt_made, kept commitment, and so on. Each signal points back to the narration it came from, so the chain from observation to claim is always visible.

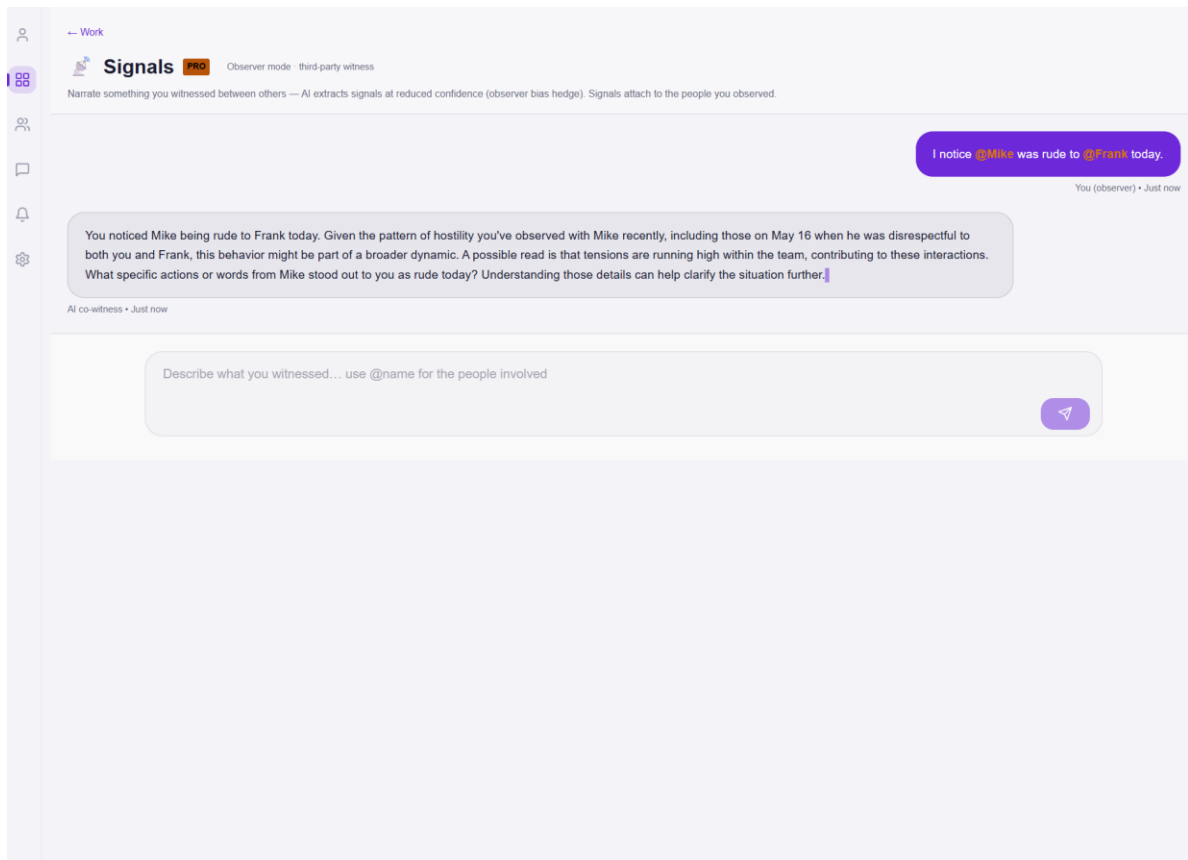


Figure K, Signals is similar to observation analysis noticing entities interactions between two different people that are already assigned to the space interaction. This is similar to notice Bob talk to Alice and her reaction seemed interesting enough to notate and communicate to the AI, so long as the user invoked the user with @ button than those signals will be appended to the user profile increase the model.

Mirror - the AI learns by asking (Subscription Only)

Mirror is the surface where the AI asks the user clarifying questions to refine its understanding of the space. It is the inverse of the journal, the user listens, the AI speaks. Questions are scoped to gaps the substrate has identified, entities the system suspects but does not yet have evidence for, ambiguities in earlier narrations, places where signals could be interpreted multiple ways. Mirror is how the substrate gets sharper without the user having to constantly self-correct.

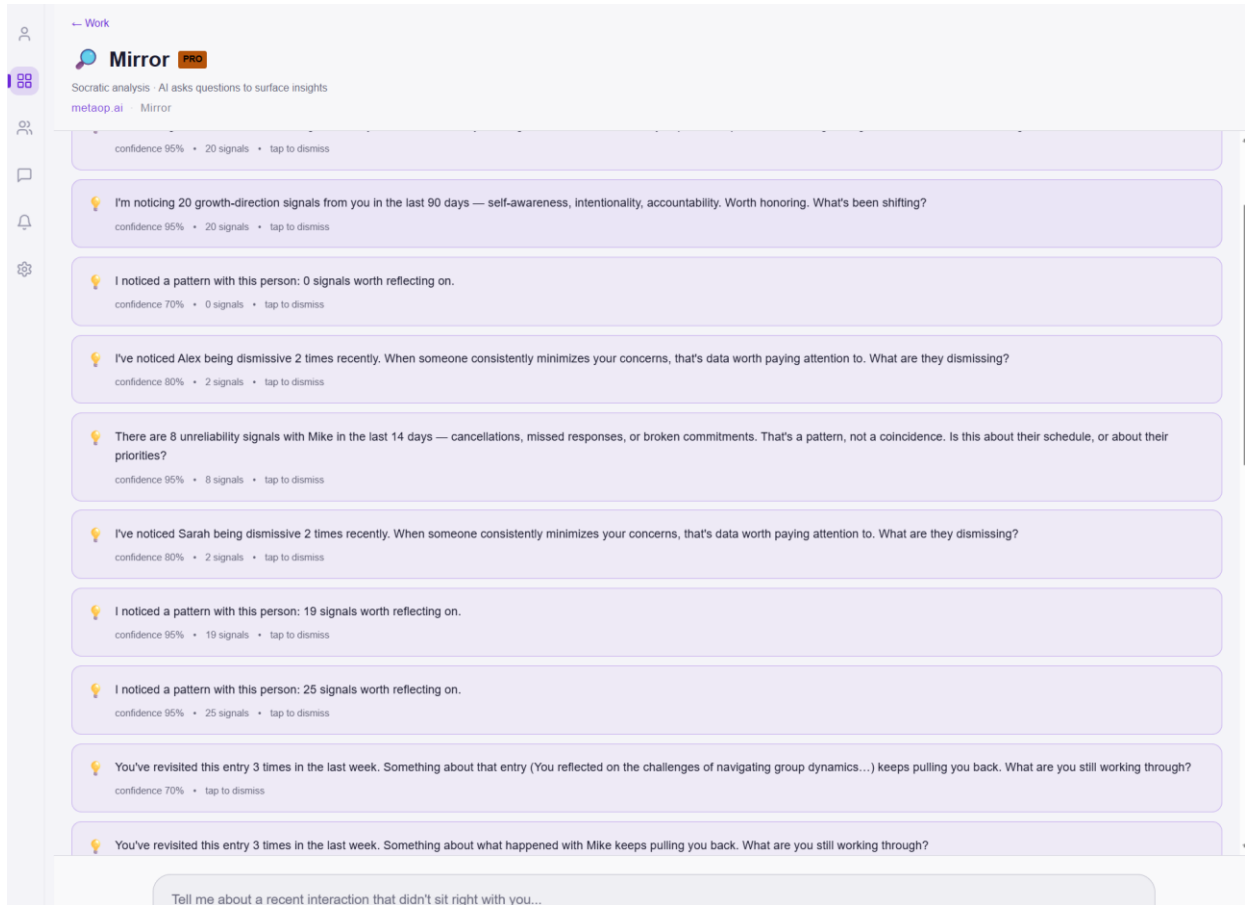


Figure L, Mirror is a paid-tier feature where the AI asks Socratic questions, grounded in a space's signal history and the user's emotional patterns, to refine its model of the user (emotional baseline, reactivity, self-awareness) within that space.

360 - the balanced view (Subscription Only)

360 is the surface that holds the adverse and the positive side by side. We deliberately do not let the system surface only what is wrong. Outcomes are modeled using game-theoretic primitives, when a pattern fires, the system considers what behaviors would compound it versus interrupt it, and pairs the adverse signal with a contextual positive counterweight. The point is honesty without spiral. The substrate has seen what is happening; the user deserves both sides.

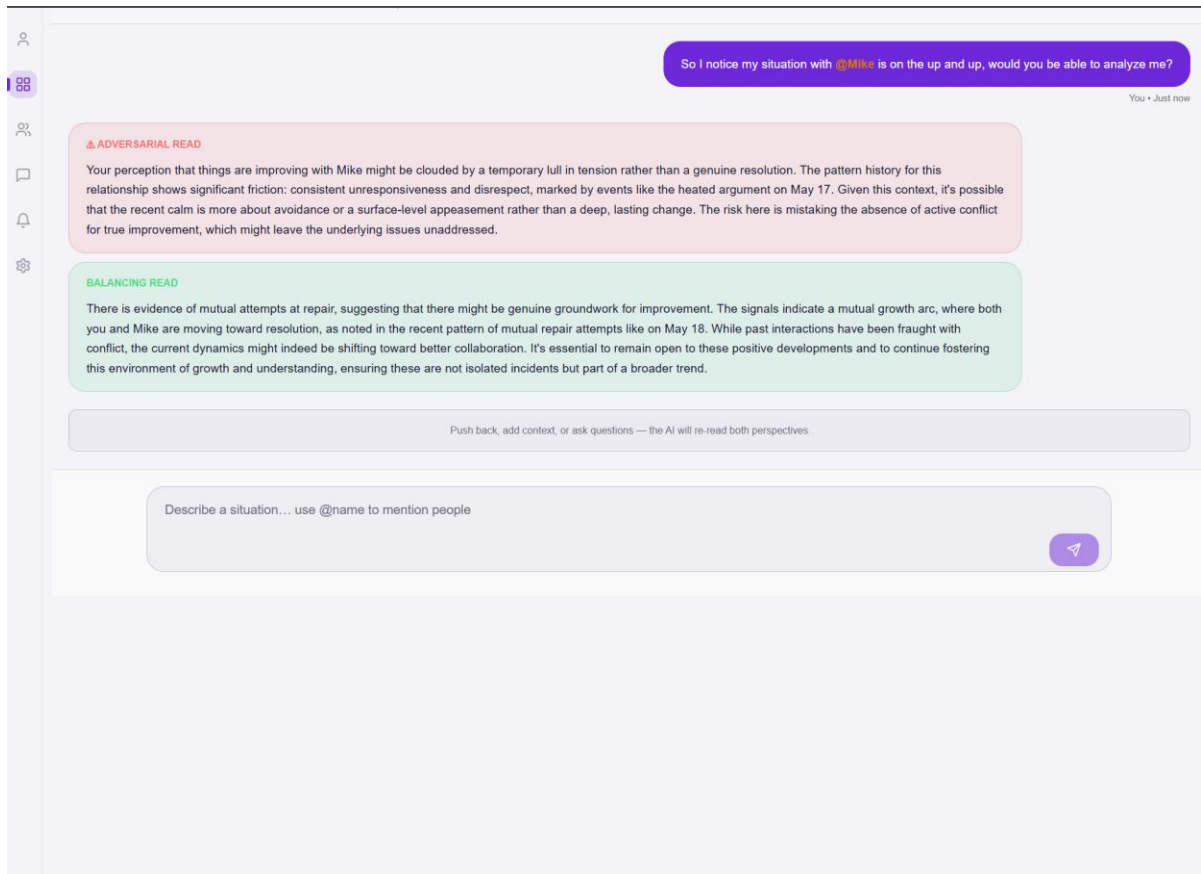


Figure M, 360 Adverse the AI will be honest with the user and be able to provide hard reality from its observations from the users' own narration, however we preface this feature with a warning that the user must agree to and we also accommodate the good things the users does do.

Outcomes (Subscription Only)

Entities, Users, and Spaces function as evolving profiles modeled over time using signal intelligence extracted from the user's own narrations. The system is intentionally perspective-aware: it primarily models the user's interpretation of events rather than claiming objective truth. One of the more interesting analytical capabilities is probabilistic outcome modeling. While the system does not provide verdicts or definitive judgments, both for ethical reasons and because no AI system is qualified to fully arbitrate human reality, it can simulate likely outcomes based on hypothetical scenarios and the continuity patterns already present in the substrate.

Initially, the system generated extremely large outcome trees with hundreds or even thousands of possible trajectories, but this became cognitively overwhelming and operationally noisy. The output was therefore constrained to the top few most probable outcomes based on the current continuity state and pattern density.

This capability is useful because it creates a feedback mechanism between the user's internal worldview and observed continuity. If the likely outcomes generated by the system consistently diverge from what the user expects, it may suggest that the user's narration, assumptions, or interpretation of events are becoming skewed relative to the broader behavioral patterns present in the substrate. The system therefore does not attempt to declare truth, but it can help surface possible misalignment between perception and continuity over time.

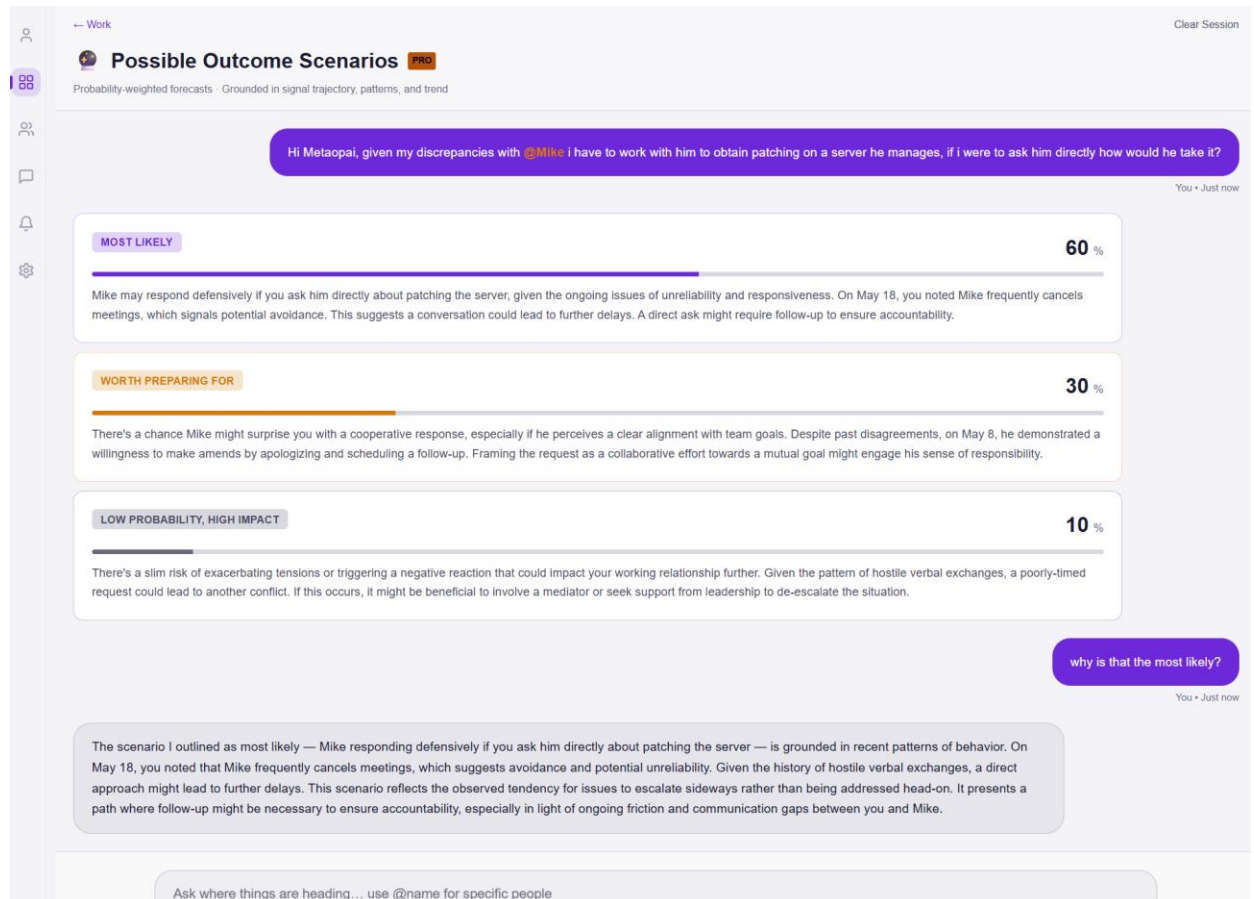


Figure N, Outcomes, because the system models the user's narrated people, spaces, and recurring patterns, it can simulate likely outcome branches for a proposed event. The UI surfaces only the top few plausible branches with confidence and uncertainty, rather than claiming exhaustive prediction.

Space Intel Report

The Intel Report is the system's narrative analytical summary of the space, group dynamics, organizational stress, social structure, environmental drift, multi-entity climate. It is the only place in the product where we use a Tavily-backed RAG enrichment by default, third-party public context that can sharpen the system's read of, say, an industry

or a public organization the space refers to. Enrichment runs only when the user is on a surface that has explicitly opted in, and the enrichment status is shown when it failed (so the user always knows whether they are reading enriched or substrate-only output).

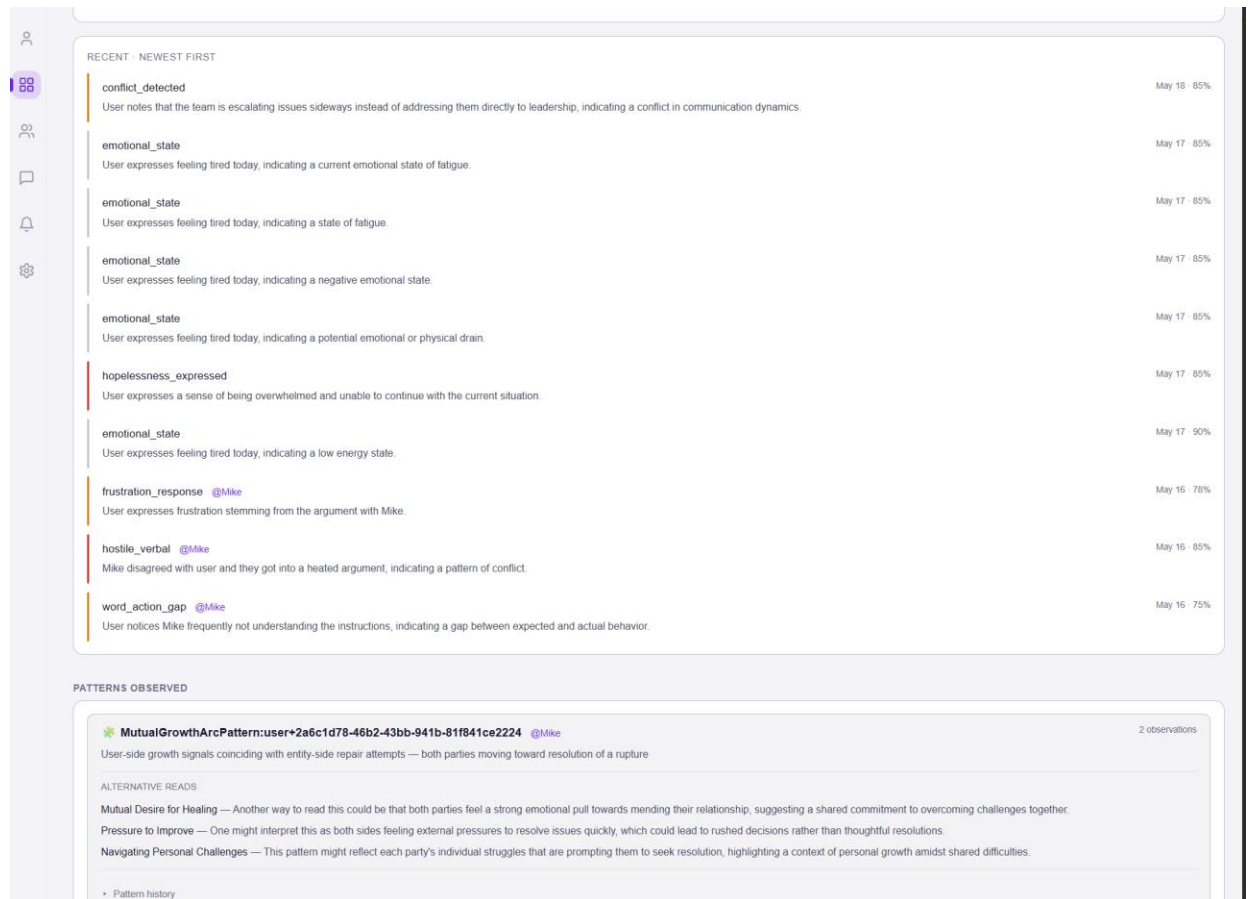


Figure O, Space Intel Report is extensive with data and metrics in the image above shows the signals and patterns observed and appended data of when.

Asking questions without changing the substrate

The Space Analytics surfaces are read-only over the substrate. The user can ask analytical questions of a space without those questions writing new patterns or signals back. This is a deliberate boundary. The journal is the right path; the analytics are the read path. Mixing them would create a recursive contamination problem where the system's analysis of itself becomes evidence for its analysis of itself. We refused to do that.

Entities - the dyadic surface

Every entity has its own page. The page is the read view: an Entity Intel Brief at the top, a Detected Patterns list with confidence and evidence, a Signal Trend chart, and the option to

open the Entity Chat. The Entity Chat (entity journal) is a dyadic conversation about the relationship between the user and that entity, scoped at ENTITY rather than SPACE. It is the right surface when the user wants to think about one person specifically, not the environment they share.

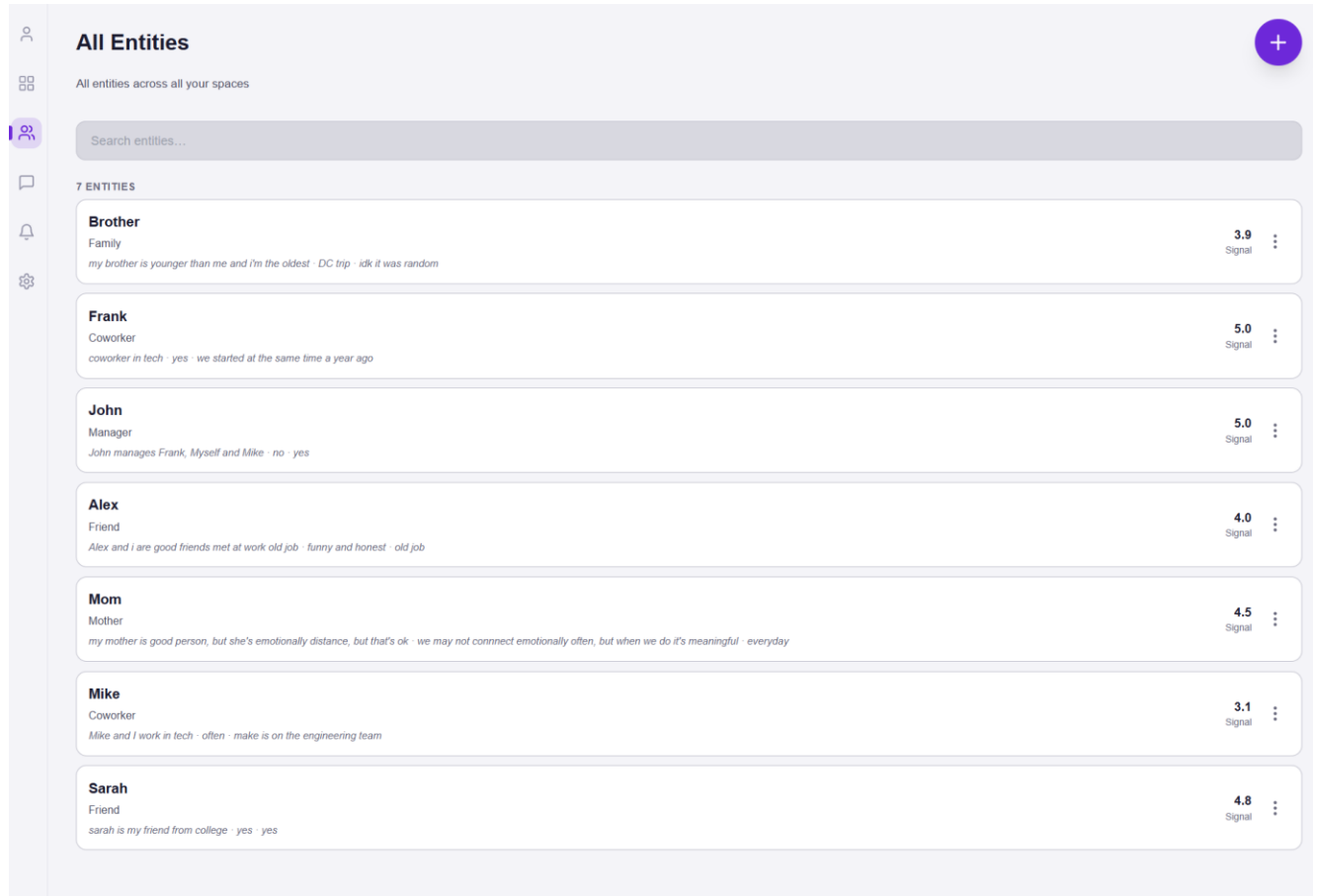


Figure P, Entities listed contains a list of the people that the user can model think of them as empty profiles and as the user journal within a space or the entity journal it will extract signals.

Entity Intel Reports + Tavily enrichment

Entity Intel Briefs draw on the same Tavily-backed enrichment as Space Intel Reports, with one strict guardrail: the system never speculates about an entity beyond what the user has narrated or what is verifiably public. When enrichment fails or returns no relevant material, the brief says so. We would rather show a gap than fill it with confabulation.

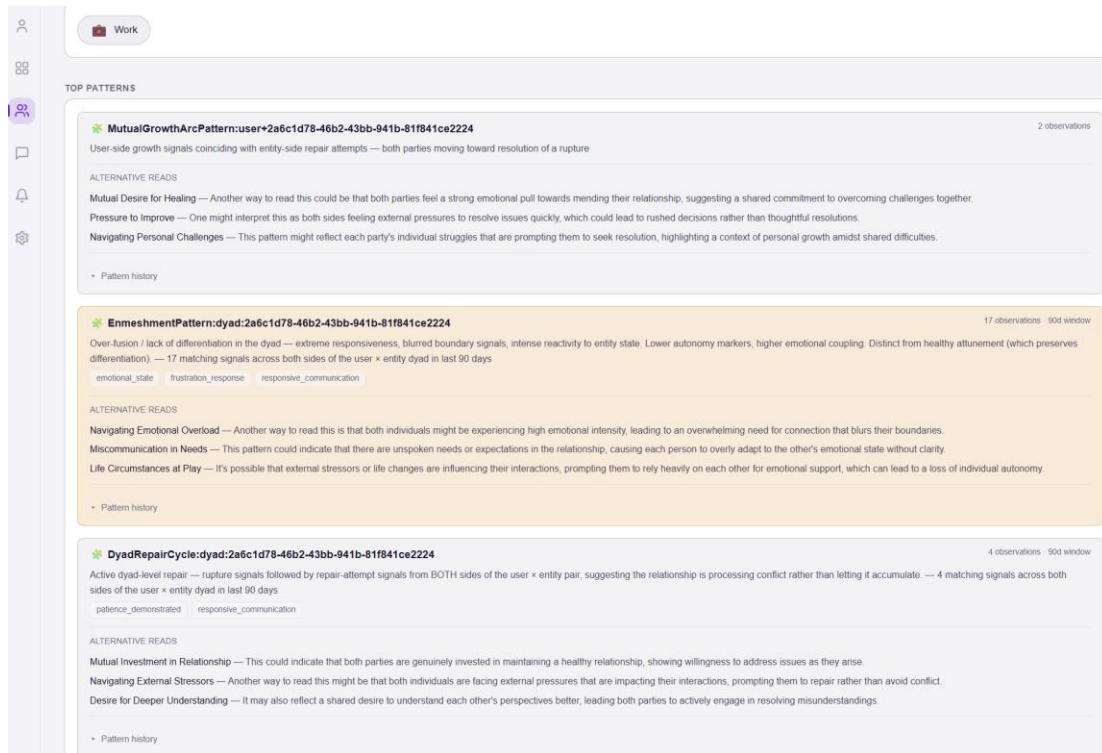


Figure Q, Per Entities intel report shows relationship signals and patterns.

Activity - the cross-cutting view

Activity is one pane that shows every signal, pattern, compound pattern, multi-pattern, recurring pattern, and trend the system has observed for the user, annotated with which entity or space the observation came from. It is the surface a user opens when they want to see what the system has seen, without committing to a specific space or entity. Items are filterable by scope, by recency, by severity, and by entity. Multi-space patterns surface as a single row with a small dropdown listing the spaces involved (opaque background, hover-only behavior, no pinning). The view is read-only, interacting with an item navigates to the appropriate detail page; nothing here writes back.

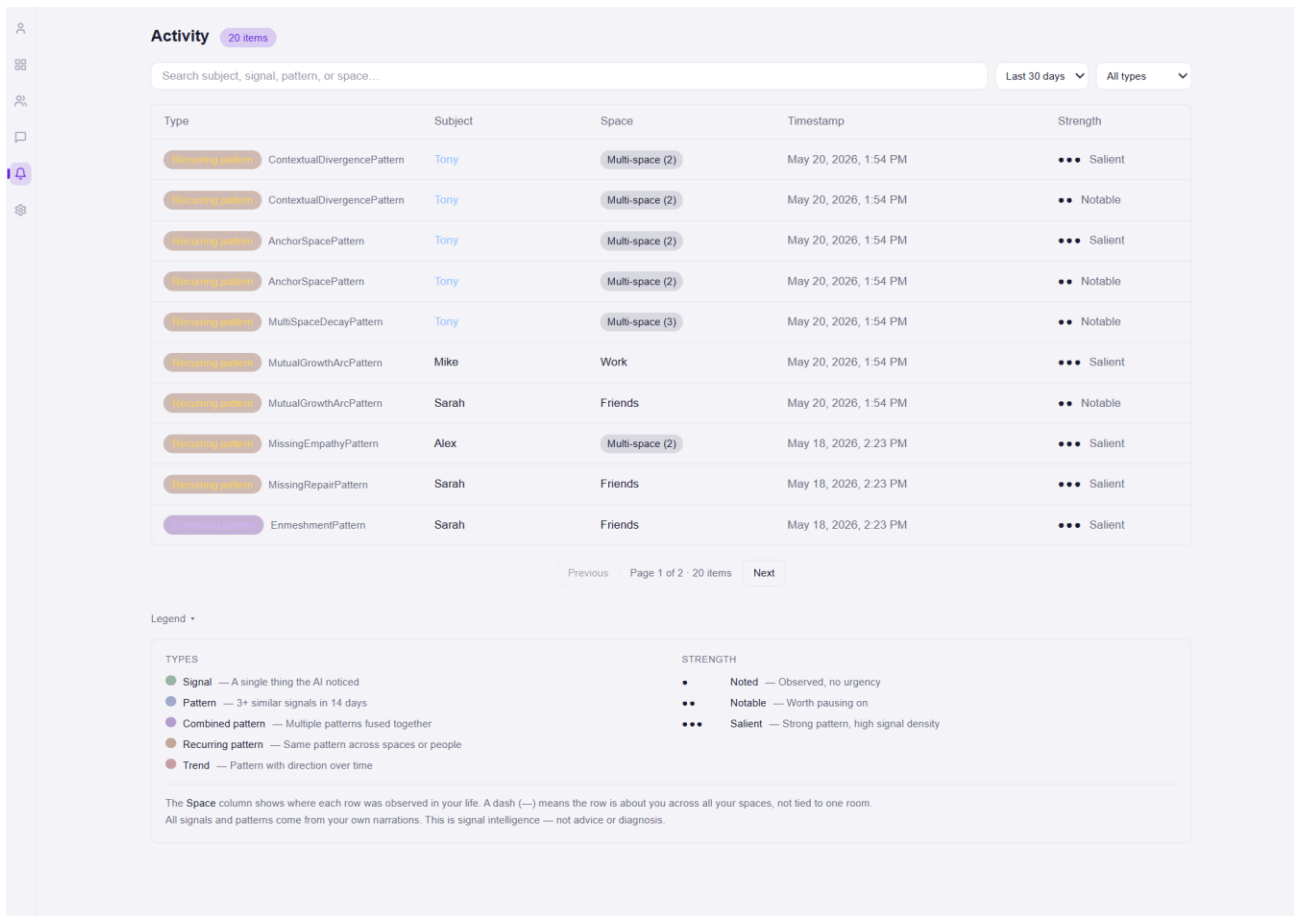


Figure R, Shows all signals, patterns, multi patterns, compounded patterns across the users, entities, and spaces in one pane.

Cross-Space - the whiteboard

Cross-Space is the analytical surface of MetaOpAI. It functions like a whiteboard where users can pull together multiple journals, entities, and environments to explore correlations and patterns across them with the help of AI. Nothing inside Cross-Space is permanent memory, it is a temporary analysis surface designed to help users make sense of the signals and continuity already present in the substrate.

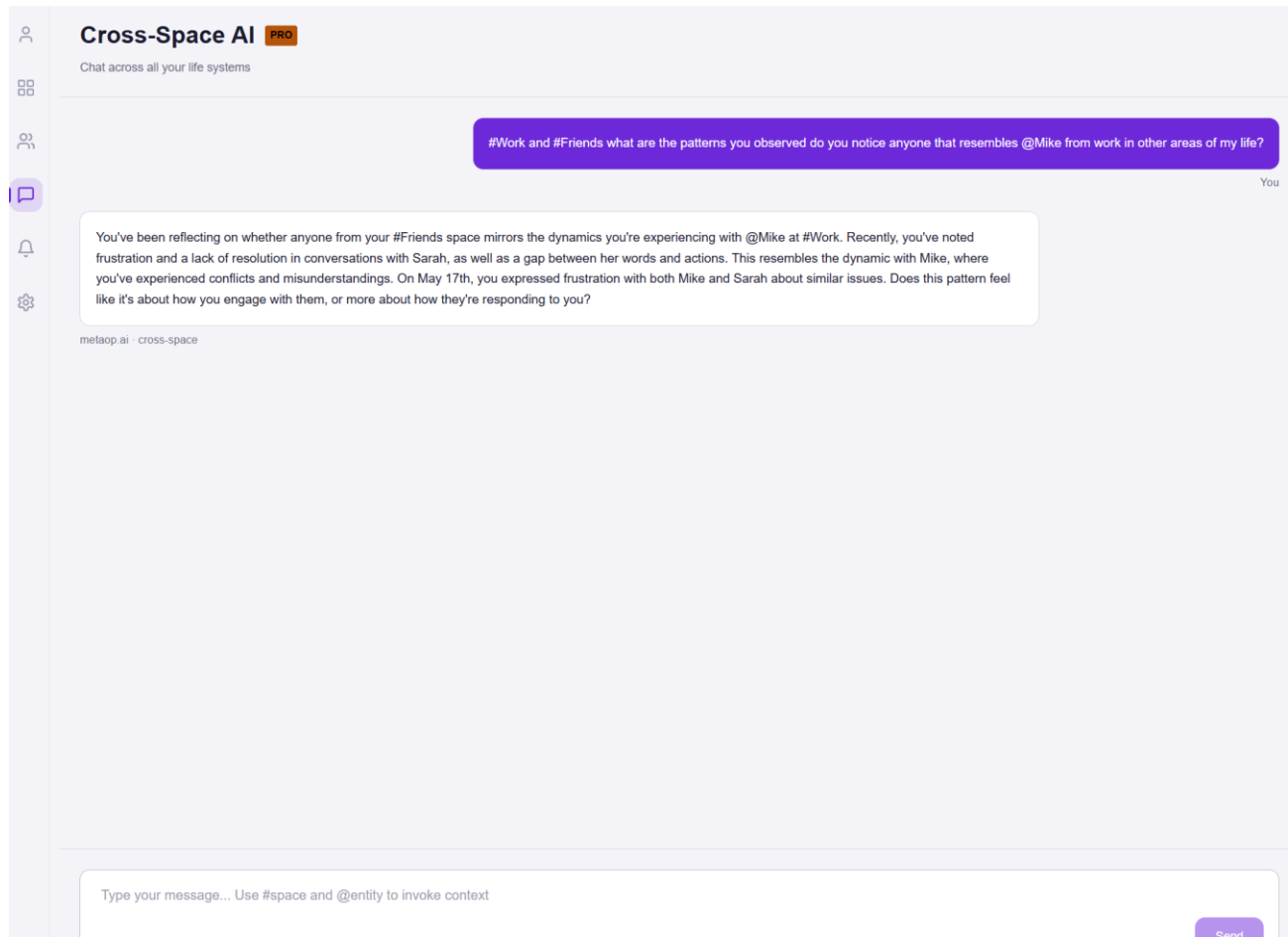


Figure S, Whiteboard analytics where the user can import journals and entities and run queries and correlate across multiple instances of their lives to understand patterns.

Cross-Space is the analytic surface for ad-hoc correlation. The user imports one or more spaces and entities (the mental model is like Splunk indexes or data sources), then runs an analytic conversation that draws on the union of the selected scopes. Crucially, the whiteboard is ephemeral. It is not saved as a journal. Nothing the user writes here lands in the durable substrate. We built it this way deliberately, there are questions the user wants to ask of their own continuity without those questions becoming part of their continuity. The whiteboard is the surface for those questions.

Settings - what the user controls

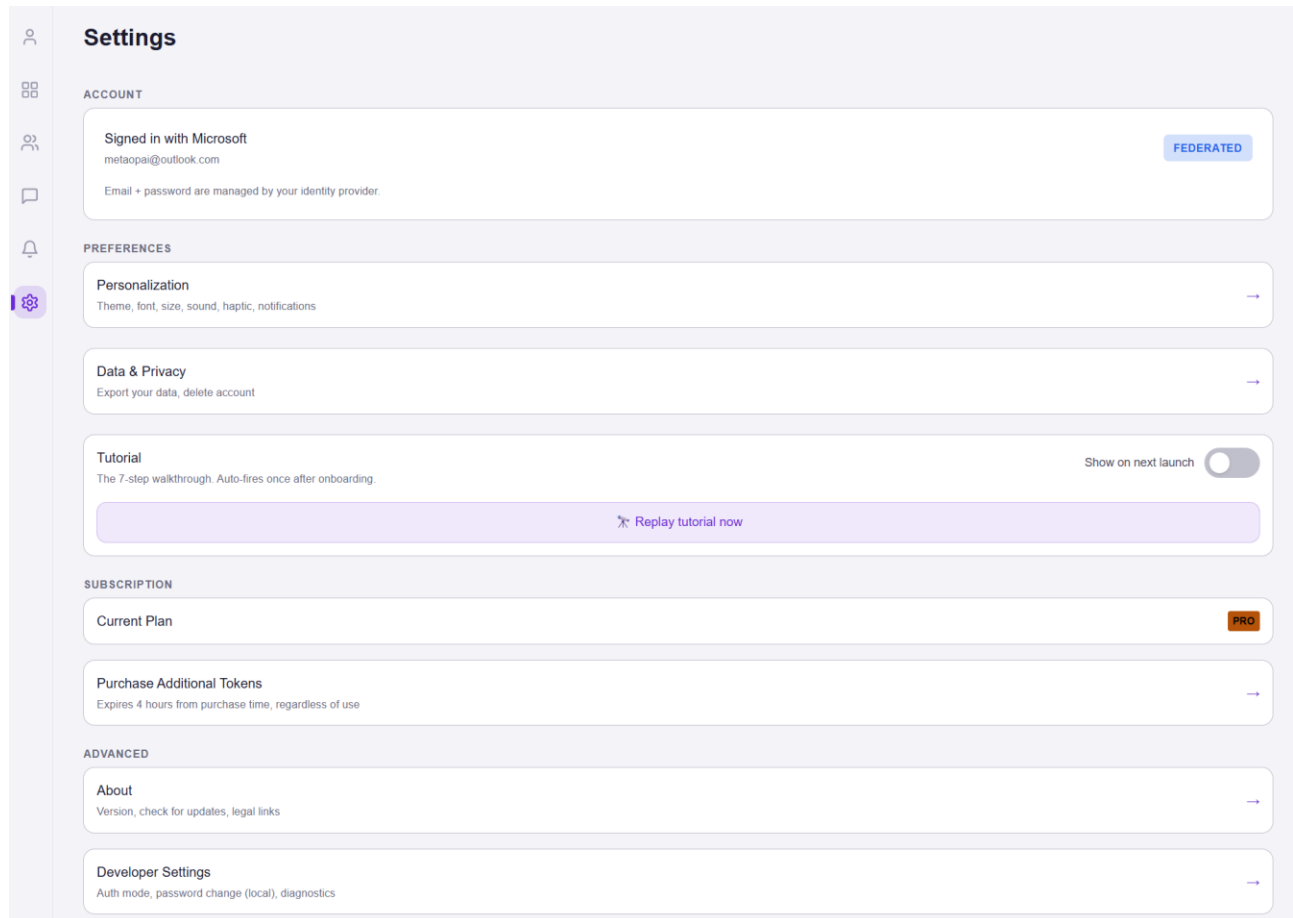


Figure T, Settings Pane

Settings is short on purpose. The major dials:

- Account, federated auth only in production: Google, Apple, Microsoft. We do not store passwords. The dual-mode auth is environment-gated; the local-password mode exists only for development.
- Subscription, the user's tier (Free / Entry / Mid / Pro), usage gauge against the daily cap, the option to change tier with the archive-keep-N workflow for resources over the new tier's caps, and the option to purchase a Burst when the cap is hit.
- Data Privacy, Users can export their complete substrate data as structured JSON, disable Behavioral Intelligence collection to stop new behavioral metrics from being generated, and request full account deletion. Existing cognition records remain intact unless explicitly deleted by the user. These controls are designed to support core GDPR principles including data portability, user control, and the right to erasure.
- Notifications, surface preferences, daily summary opt-in, crisis-resource visibility.

- Legal, Terms of Service, Privacy Policy, Ethics Statement, Crisis Resources. These are not buried in a footer. They live in Settings because they are part of how the product works.

Behavioral Capture - Honestly

Why we collect typing rhythm, latency, and entry texture, and the controls users have over it.

Two narrations with the same words can mean very different things. A sentence typed quickly with high error rate at 1 a.m. is not the same data point as the same sentence typed slowly at 10 a.m. If the substrate is going to ground its cognition in what the user actually said, it has to ground it in how the user said it as well. That is what the behavioral capture model is for.

Data & Privacy

Export everything we have on you, or delete your account entirely.

EXPORT

Download everything
Your profile, signals, journal entries, estimators, and settings — packaged as one or more .zip files. Large datasets are split into batches so each download stays fast and resumable.

[Estimate my export](#)

BEHAVIORAL ANALYTICS

Capture typing patterns

Records typing speed, response time, message length, and paste detection on each entry. Powers the Your Rhythm and Growth Signal patterns on your profile.

When OFF: capture stops immediately (no future data collected). Your Rhythm and Growth Signal sections are hidden from your profile until re-enabled. Past data is preserved — use Export above for a copy, or delete your account below to remove all past behavioral data.

CONSENT HISTORY

A record of every time you accepted or declined the Terms of Service or Privacy Policy. This is the audit trail used for compliance and for your own records.

✓ Accepted Apr 27, 2026, 07:06 PM

Terms v1.0 · Privacy v1.0 · via microsoft · from 172.19.0.1

Records are append-only and cannot be edited or deleted. They are retained for compliance purposes for up to 7 years even after account deletion.

DELETE ACCOUNT

Permanently delete your account
Removes everything — profile, journal entries, signals, estimators, billing history. Cannot be undone. Backups rotate within 90 days. Export your data first if you want a copy.

[I want to delete my account](#)

Figure U, Data Privacy, where a user can export their data, disable behavioral metrics from collecting, and they can delete the account.

What we capture

- Typing rhythm, the cadence of keystrokes, including pauses and bursts.
- Response latency, how long the user spent between the AI's last message and their reply.
- Sloppiness drift, error rate over a session, normalized against the user's baseline.
- Stress state, a composite signal derived from latency, error rate, and word-choice markers.
- Word choices, sentiment-loaded or schema-relevant vocabulary that the extractors weight.
- Entry length trend, average narration length over time as a proxy for engagement and energy.

How we use it

Behavioral data informs the substrate's reading of context. A signal extracted from a high-stress narration carries different weight than the same signal extracted under normal baseline conditions. Patterns weight behavioral metadata when computing confidence. The composer uses behavioral context to choose the appropriate tone, terser when the user is short on energy, more reflective when the rhythm suggests they want to think.

How we govern it

Capture is gated by a single toggle in Settings, Disable Behavioral Intelligence. When the toggle is off, no behavioral metrics are written for the session; the metric columns are NULL. We do not show the toggle as a display gate (where the data still gets collected and the user just cannot see it). It is a collection gate. Disabling behavioral intelligence reduces the substrate's contextual sharpness on those turns; we say that plainly. We do not pretend the feature is free.

Per-category deletion of behavioral data was explicitly declined as a design decision, the metric is meaningful only relative to baseline, and selective deletion would corrupt the baseline. Users can delete their entire account or export everything to JSON; we do not offer a middle path here because the middle path would create silently broken cognition.

PART III

Integrated Product Additions: User Mental Model and Premium Wedge

This section folds the product-whitepaper lens into the unified paper. The product is not merely a journal with an assistant attached to it. The interface is the ontology made visible: Profile for USER cognition, Entities for ENTITY cognition, Spaces for SPACE cognition, and Activity / Cross-Space / Evidence workflows for RELATIONSHIP_PAIR cognition. The user does not need to understand the KRL to benefit from it; they experience it as surfaces that answer concrete human questions.

The user mental model

Users naturally organize life into three domains: themselves, the people around them, and the environments they operate inside. MetaOpAI maps those intuitive domains into durable substrate scopes. This is why the product feels familiar even though the architecture underneath is technical: the database model follows the way users already narrate reality.

When the user asks, "What is going on with me?", the Profile and Growth Journal surfaces read and write USER-scope cognition.

When the user asks, "What is happening with this person?", the Entity surface reads ENTITY and RELATIONSHIP_PAIR evidence without turning the person into a diagnosis.

When the user asks, "What is happening in this environment?", the Space surface reads systemic signals, events, and meta-context about that environment.

When the user asks, "Where else is this pattern appearing?", Activity, Cross-Space, Evidence Chain, and KRL views expose the relationship between scopes.

Premium product wedge

The premium wedge is not a single feature. It is compounding continuity. Free or lightweight usage can demonstrate journaling, basic reflection, and early signal capture. The paid product becomes more valuable as the user accumulates enough structured evidence for cross-scope correlation, relationship-state movement, timeline reconstruction, pattern confidence, and auditability.

KRL Explorer turns the substrate into a readable library of surfaced cognition.

Graph Explorer turns the same substrate into a node-based map of people, spaces, relationships, and evidence density.

KRL Formation shows which records compounded into a pattern rather than asking the user to trust a black-box answer.

Evidence Chain gives the receipts: the actual user-authored source moments that produced the insight.

Cross-Space is the power-user surface: it asks what is repeating across environments, people, and relationship pairs.

Product boundaries

MetaOpAI should remain clear about what it is and is not. It is not therapy, diagnosis, a medical device, or a replacement for human judgment. It is a signal-intelligence product: it names observable interaction patterns, shows evidence, uses confidence, tracks contradictions, lets the user correct the model, and avoids fabricating meaning when the substrate is sparse.

That boundary is product-defining. The system can say, "this interaction pattern appears to be repeating," but it should not say, "this person is a narcissist." It can show evidence of avoidance, asymmetry, withdrawal, repair attempts, or repeated conflict. It should not convert those observations into identity verdicts about a person.

Architecture

How the ontology shows up in the UI, onboarding, journals, analytics, the surfaces a user actually touches.

Architecture Overview

Seven layers, each separately versionable. The substrate is the moat; the model is the rented edge.

Layer [1]: How data gets into MetaOpAI (Growth Journal, Entity Journal, Space Journal)

The space journal is the richest ingestion surface because a space contains the user, the people inside it, the relationships between them, and the environment defines the constraints. A single narration inside a Work space, “Mike cut me off again and I shut down”, can simultaneously generate a USER signal about avoidance, an ENTITY signal about Mike, a RELATIONSHIP_PAIR observation about the interaction, and a SPACE-level signal about the environment itself. One narration can therefore write across multiple cognition scopes at once.

On the read side, however, the journals remain deliberately separate because each one represents a different cognition identity rather than simply a filtered view of the same data. The growth journal speaks introspectively about the user. The entity journal speaks dyadically about the relationship between the user and another person. The space journal speaks systemically about the broader environment and the people inside it. All three journals’ types share the same substrate underneath, but they are intentionally designed to think and speak differently. Collapsing them into a single conversational surface would flatten those distinctions and force every interaction into the same systemic voice, weakening the introspective and relational perspectives that make the cognition architecture useful in practice.

How Extractors Work

The system performs five LLM calls per conversational turn: one synchronous composer call that generates the user-facing reply, and four asynchronous extractor calls that run in the background after the response has streamed. These extractors are intentionally organized by record type rather than by subject dimension or scope.

The extractors consist of:

- Signal Extractor: which identifies behavioral markers such as withdrawal, avoidance, or kept commitments;
- Event Extractor: which identifies discrete moments and occurrences such as arguments, meetings, or transitions;
- Meta-context Extractor: which captures the user’s framing and interpretation of events, separating objective context from subjective narrative;

- Relationship-pair Extractor: which analyzes the dynamics between the user and a specific entity, such as trust erosion, reciprocity imbalance, or repair attempts.

A glass with 4 oz of water is an objective fact, that is Context. But when one person calls it “half-full” and the other calls it “half-empty,” those are interpretations of a fact, that is Meta-Context. The system stores both separately: the factual state of the glass and the way people frame or perceive it. This separation is important because trustworthy cognition requires distinguishing what objectively happened from how it was emotionally or psychologically interpreted. Once these records are stored with full provenance metadata, including extractor source, version, timestamp, message origin, and confidence, the pattern engine evaluates whether enough evidence exists to create, reinforce, decay, or expire higher-level patterns within the substrate.

Importantly, the extractors themselves do not determine whether a record belongs to the USER, ENTITY, SPACE, or RELATIONSHIP_PAIR scope. Instead, the model emits semantic references such as “user” or “Mike,” and a downstream resolver maps those references onto known subjects within the substrate before assigning scope during persistence. This creates a grid-like architecture in which extractors function as the analytical lens while scopes function as the subject classification layer:

EXTRACTOR	USER	ENTITY	SPACE	RELATIONSHIP_PAIR
Signal	✓	✓	✓	
Event	✓	✓	✓	
Meta-context	✓	✓	✓	
Relationship-pair				✓

This separation keeps extraction logic orthogonal to subject organization, allowing the same extractor to produce records across multiple scopes depending on the narration being analyzed. All extractor outputs are forced through strict structured-output schemas

with typed enums, required fields, and validation enforcement. Invalid outputs are retried once with validation feedback and discarded entirely if they fail again, reflecting a core design principle of preferring omission over substrate corruption. Before persistence, records are filtered through centralized confidence thresholds, unresolved entities are rejected rather than auto-created, and provenance metadata is attached as a mandatory invariant, including extractor type, extractor version, source message, and write timestamp. Once persisted, the pattern engine evaluates whether newly ingested material crosses density and temporal thresholds sufficient to create, reinforce, decay, or expire patterns, after which a settle marker is written so downstream cognition and retrieval systems know the substrate has reached a coherent post-turn state.

Layer [2]: Knowledge Representation Layer (KRL)

KRL is the continuity substrate. Four subject scopes (USER, ENTITY, SPACE, RELATIONSHIP_PAIR) by five-layer kinds (Signal, Event, Meta-Context, Context, Pattern). Twenty addressable cells per user. Every record carries provenance, extractor name and version, source message id, timestamp, confidence. The matrix shape is what lets the orchestrator ask precise questions of the substrate. We do not blob everything into a single retrieval index and hope. We ask: what does the system know at the dyad level between this user and this entity, and only retrieve that.

The substrate is the part we are proudest of, and the part hardest to see. This chapter is the long version.

What the KRL is solving

A matrix. Four subject scopes by five-layer kinds, twenty addressable cells per user. The subject scopes answer a simple question: **who or what is this about?** USER represents the person themselves. ENTITY represents another person in the user's life. SPACE represents an environment such as Work, Family, or a friend group. RELATIONSHIP_PAIR represents the dyad between the user and another entity, resolved through a deterministic relationship identifier so the same pair always maps back to the same continuity record.

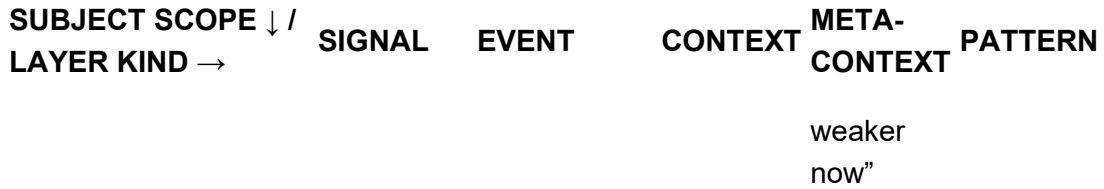
The layer kinds answer a different question: **what kind of knowledge is this?** Signals are lightweight behavioral observations extracted from narration, such as withdrawal, avoidance, or repair attempts. Events are discrete narrated moments, arguments, conversations, meetings, transitions, or milestones. Context represents stable facts that are expected to remain relatively durable over time, such as "Sarah is your sister" or "Mike is your manager." Meta-Context stores interpretive framing rather than objective

reality, for example, “I think Sarah is becoming distant.” The system deliberately separates factual events from emotional interpretation because interpretations can drift, evolve, or later prove incorrect. Patterns sit at the highest layer: they are aggregate abstractions computed by the pattern engine across multiple signals, events, and time windows.

The result is a structured cognition matrix rather than a conversational blob. Instead of storing raw chats and replaying them indefinitely, the substrate converts narration into typed records that can be selectively retrieved, weighted, decayed, correlated, and reasoned over independently. Every record exists inside a precise location within the matrix: a specific subject scope and a specific layer kind. This allows the orchestrator to ask targeted questions of the substrate such as: “What recent dyadic patterns exist between this user and this entity?” rather than retrieving an entire conversation transcript and hoping the language model notices the important parts.

Concretely, each extracted record is stored as a structured row rather than a free-floating string. A record contains the subject scope, subject identifier, layer kind, optional sub context, temporal scope, confidence score, polarity, perspective, structured JSON payload, source message reference, evidence count, contradiction references, provenance metadata, timestamps, and lifecycle state. Every field exists because the cognition system eventually needed to query, weight, audit, decay, correlate, or reason over that dimension of continuity. The philosophy is deliberate: conversations are transient, but structured cognition compounds.

SUBJECT SCOPE ↓ / LAYER KIND →	SIGNAL	EVENT	CONTEXT	META- CONTEXT	PATTERN
USER	“I shut down”	“Panic attack Tuesday”	“You work in finance”	“I feel stuck lately”	Withdrawal Spiral
ENTITY	“Mike interrupted”	“Argument at lunch”	“Mike is your manager”	“Mike seems distant”	Reciprocity Drift
SPACE	“Team tension rising”	“Layoff meeting”	“This is your Work space”	“The culture feels unstable”	Organizational Stress
RELATIONSHIP_PAIR	“Avoidance between you two”	“Repair conversation”	“User Mike dyad”	“Trust feels”	Trust Rebuilding



This matrix structure is the foundation of the cognition substrate. It allows the system to preserve continuity without endlessly replaying raw conversations back into the language model, turning narration into structured, queryable intelligence rather than accumulating transcript history.

How we solve it – the design choices matter

Deterministic addressing instead of similarity search. Because every record carries a subject scope, subject identifier, and layer type, the orchestrator can ask the substrate precise questions like: “show me recent Signal-layer records about the relationship between this user and Mike from the last 30 days above a confidence threshold.” This is a direct indexed lookup, not a broad semantic search over conversation embeddings. The system does not dump everything into a single vector index and hope the right memories appear; it retrieves exactly the continuity cell it needs.

Some layers are designed to evolve over time. Context and Meta-Context can drift, and that drift is itself meaningful. A user’s interpretation of someone may gradually shift from “distant” to “hostile” to “malicious,” and the progression matters just as much as the current state. To preserve that continuity, previous versions of changing records are stored in a history table before updates occur. Signals and Events do not require the same treatment, signals naturally decay over time, while events remain append-only historical facts, so the system only versions the layers where longitudinal drift is cognitively important.

Contradiction is treated as information rather than an error. The substrate can hold multiple competing interpretations or conflicting records simultaneously instead of overwriting one with another. If a user’s industry is described as “banking” in one narration and “telecommunications” in another, the system preserves both and allows the composer to surface the inconsistency directly. This lets the cognition layer ask clarifying questions rather than silently collapsing ambiguity into false certainty.

Every record also carries strict provenance metadata. The system records which extractor produced the record, which extractor version was used, which message created it, and when it was written. No row can exist without provenance, and automated QA checks fail if that invariant is violated. This creates a fully traceable cognition substrate where downstream patterns and observations can always be linked back to the original narration that produced them. Synthetic bootstrap cognition used during onboarding is intentionally stored outside the durable substrate so temporary scaffolding can never silently contaminate long-term continuity.

Finally, the substrate favors soft deletion over destruction. Records are typically hidden or deprioritized rather than permanently erased. Older continuity may become relevant again later, so the system treats the past as something to be weighted down over time rather than discarded entirely.

Why is it useful

Because it makes scope-precise cognition possible. When you open the entity journal for Mike, the system retrieves the ENTITY and RELATIONSHIP_PAIR cells for Mike and *nothing else*, not your work-space dynamics, not your self-reflection, not the eleven other people you've mentioned. The prompt that reaches the model is small, dense, and exactly on-subject. That is what keeps the window clean, the cost flat, and the signal-to-noise high. The matrix shape is not bookkeeping; it is the thing that lets every other layer be selective.

Is it enough?

No, and we want to be honest about that. The KRL is the substrate, not the whole product. By itself it is an organized archive. It becomes cognition only when the Pattern Engine reads it to recognize recurring dynamics and the Orchestrator reads it to assemble scope-correct context. The KRL is necessary but not sufficient. It is the foundation the rest of the system stands on, and a foundation is not a house.

There are also open edges. At very large per-user scale the retrieval shapes will need revisiting. We have validated the structure against real and synthetic traffic; we have not yet seen it under thousands of users with years of accumulated history.

How it performs

The common query paths are backed by three composite indexes: (user, subject kind, subject_id, layer) for list-by-subject, (user, layer, created_at) for recent-by-layer scans, and (user, subject_kind, subject_id, subcontext) for single-field drift queries. At current scale the aggregate read path is sub-millisecond, when we benchmarked porting the Python grouping to SQL GROUP BY, the Python version actually won at the test scale ($\approx 0.096\text{ms}$ vs $\approx 1.15\text{ms}$ at 449 rows), because the per-row cost of the SQL CASE-and-regex work dominated. We kept the Python path and documented the crossover point

where SQL would win (hypothesized around 5,000 rows) to revisit on Pro-tier telemetry. The point is not that we are fast everywhere; it is that we measure before we optimize and we know where the crossover lives.

The KRL is our answer to both failure modes at once. It is a structured, addressable substrate that knows *who or what* every piece of continuity is about, *what kind* of thing it is, *how confident* we are in it, *where it came from*, and *how it has changed over time*. It is the difference between a shoebox of photos and a labeled archive

Is it domain specific

The KRL substrate is domain-general. The same matrix structure (four scopes × five layer types) can support many different areas, personal relationships, team dynamics, patient care, research projects, or others, without any changes to the core system. Only the pattern templates and parts of the extractor vocabulary are currently tuned for relationships and self-reflection. This clean separation is deliberate: it lets the substrate remain general-purpose while allowing MetaOpAI to expand into new domains or become a platform in the future, without rebuilding the foundation.

Layer [3]: Pattern Engine

The pattern engine is the layer responsible for turning raw signals and events into recognizable, higher-order dynamics. The pattern engine continuously scans the KRL and produces typed pattern artifacts when aggregate conditions are met. Each pattern template defines:

- Which signals and events it monitors
- Minimum density and confidence thresholds
- Relevant time windows
- Awareness tier (how prominently it should surface)
- Action tier (how strongly it influences composer behavior)
- Lineage tracking back to supporting evidence

Patterns fire only when evidence crosses a deliberate threshold — never on isolated events. They decay gracefully over time using severity-weighted time-to-live logic. The engine explicitly supports contradictory patterns co-existing (e.g., signs of both trust rebuilding and withdrawal). Resolving or presenting those contradictions is left to higher layers. This disciplined approach is

what allows MetaOpAI to surface meaningful longitudinal insights — trust trajectories, communication cycles, emotional state transitions — without overclaiming or hallucinating patterns.

Pattern Engine

Confidence-bounded, evidence-gated, lineage-tracked. Patterns are the abstraction layer, not the output.

A pattern, in our usage, is a typed artifact that fires when an aggregate condition crosses a threshold across signals that span days or weeks. We do not declare patterns from single utterances. That is what the earlier generation of mood-detection systems did, and it produced exactly the noise you would expect.

Patterns are stratified by tier. Tier 1 are signal clusters, N signals of one type across a window. Tier 2 are canonical/relational template patterns over a single subject (entity, user, space, or relationship-pair). Tier 3 are compounds, combining Tier-2 patterns across subjects. Tier 4 are second-order trajectory patterns about how other patterns are evolving over time — Intensifying, Fading, or Reversal. Higher tiers require lower-tier patterns to already exist. The ladder enforces a discipline against premature thematization.

Confidence - Evidence, Lineage

Every pattern carries three discipline fields. The confidence score reflects how strong the evidence is and decays as old evidence ages out. The evidence list points to the signals that justify the pattern, so the user (or an engineer auditing) can trace any pattern back to the raw narration that produced it. The parent-pattern lineage records which lower-tier patterns the higher-tier pattern depends on. If a child pattern weakens, the parent's confidence weakens too. This is not just bookkeeping; it is the audit trail.

Awareness and Action Tiers

Two further classifications govern how a pattern shows up in the product. The awareness tier (critical, high, medium, low) decides whether the pattern surfaces in the UI at all. The action tier decides whether the system suggests anything actionable. We treat suggesting action as a higher-stakes commitment than naming a pattern, so the tier ladder is deliberately steep. Most patterns we detect are surfaced for awareness but not for action.

Contradiction Handling

Patterns can contradict each other. A user might be in a TrustRebuilding arc with their partner while simultaneously showing signals consistent with WithdrawalSpiral. Both are

real. Both deserve to surface. The composer is instructed to hold both at once, not collapse them. Contradiction is information.

Time-Decay and the Signal-Noise Discipline

Most cognition systems treat old signals as permanently equal to new ones, which quietly lets the past dominate through accumulation. A difficult month from two years ago can end up shaping the system's view of you forever. We handle this differently. Signals lose weight over time instead of remaining equally important forever. Recent signals matter more than older ones, and serious signals decay slower than minor ones. A withdrawal signal from this week therefore carries more weight than ten similar signals spread across several years. Nothing is truly deleted, old records remain available for history and analytics, but they gradually matter less in the pattern engine unless they become relevant again. The result is a system that respects recency without becoming forgetful: it still remembers stable long-term facts, but when analyzing what is happening lately, it prioritizes recent continuity over distant history.

Why TTL is usually wrong here

Conventional cognition systems treat older data as having permanent value: every signal ever extracted stays in the substrate forever, and the system reasons across all of it. This sounds principled but it produces a subtle pathology. Older signals begin to dominate aggregations simply through accumulation. A user who had a difficult month two years ago, but has been doing well since, ends up with a pattern engine that is permanently colored by that historical period. Recent recovery signals are diluted by the weight of historical struggle.

Our approach: TTL as a deliberate weighting instrument

We handle time differently from most cognition systems. Signals, events, and meta-context records gradually lose weight as they age instead of remaining equally important forever. Older records are not deleted, they still exist for history views and analytics, but recent signals matter more in pattern detection than older ones. For example, ten withdrawal signals from the last two weeks produce a much stronger pattern than ten similar signals spread across two years.

Operational cognition traces are handled separately. Those traces exist only for debugging and system analysis, not for long-term continuity, so they are automatically deleted after seven days. The user's actual continuity inside the KRL substrate is not deleted over time; instead, it uses weighted decay, where older continuity gradually matters less unless it becomes relevant again.

What this lets us do

The combined effect is a system that respects temporal proximity without being amnesic. It still remembers that the user's sister moved abroad three years ago (stable fact, no decay). It still remembers patterns that fired six months ago (active record, queryable). But when answering “what's happening with my brother lately,” it weights the last few weeks much more heavily than the prior year. This matches how human cognition actually works: we remember everything in principle, but the texture of “recent” feels different from the texture of “old,” and the system's outputs honor that distinction.

What this does not do

Time-decay is not a substitute for explicit user control. Users can still ask for historical patterns, surface old events, and review their long-term trajectory. The decay is a default-weighting choice in pattern aggregation, not a hard filter that erases the past. We have been explicit with users that the system favors recent signal when forming current reads, and that they can ask to widen the lens whenever they want to see further back.

Layer [4]: Session Cache

The session cache - the synthetic cognition blob.

The synthetic cognition blob is a temporary session-level cache used to solve the cold-start problem. When a new user begins using the system, the durable cognition substrate (KRL) contains little or no continuity yet, leaving the orchestrator with very little meaningful context to retrieve. Instead of allowing the language model to invent continuity to fill those gaps, the system stores lightweight synthetic onboarding context inside a temporary conversation-scoped cache. This material is explicitly marked as low-authority scaffolding rather than durable cognition.

This tension, between cold-start usefulness and recursive synthetic drift, is what we addressed internally during Phase L. Many AI systems risk creating epistemic feedback loops where model-generated interpretations are written back into memory and later re-read as if they were observed user reality. Over time, the system begins compounding its own interpretations instead of compounding the user's actual narration. The result is a subtle form of recursive contamination where synthesized cognition becomes evidence for future synthesized cognition. MetaOpAI intentionally blocks this behavior by preventing synthetic cognition from ever entering the durable substrate.

The solution was to isolate synthetic cognition into a temporary session cache, conceptually similar to a DNS cache, while keeping durable continuity physically separate inside the KRL. On each turn, the orchestrator can retrieve relevant synthetic context related to the user, an entity, a space, or a relationship pair and inject it into the prompt as explicitly low-authority subcontext under strict token limits. This allows the system to stabilize onboarding and provide lightweight continuity assistance without treating synthetic scaffolding as durable truth.

Crucially, the two systems never merge. Synthetic cognition lives only on the conversation session, while durable cognition lives inside extracted substrate records generated from the user's actual narration. This physical separation acts as a firewall against recursive substrate contamination. As the user continues using the system across their first sessions and weeks, the durable KRL gradually becomes richer and more useful. This creates what we internally call the "hourglass effect": early in a user's lifetime the system leans more heavily on temporary synthetic scaffolding because durable continuity is sparse, but as real user-generated cognition accumulates, the orchestrator increasingly prioritizes the KRL while the synthetic cache becomes progressively irrelevant. Over time, the temporary scaffolding naturally tapers away as durable continuity compounds beneath it.

Layer [5]: External Retrieval-Augmented Generation (RAG): Tavily

Tavily is used in two different ways inside the system, and the earlier whitepaper only described one of them.

The first use is inline reference lookup during live conversations. When a user mentions something concrete and public, such as a company, CVE, public figure, or external reference, the system can perform a lightweight lookup to retrieve supplementary context for that specific turn. These lookups are cached in Redis, metered against usage quotas, and fail open: if Tavily, Redis, or the lookup itself fails, the conversation simply continues without enrichment. Tavily is therefore never on the critical path of the product. This inline lookup is deliberately scoped to the Growth journal only — it is the sole inline-Tavily surface. (The Space- and Entity-level Tavily enrichment used in Intel Reports is the separate, second use described below.)

The second use is context enrichment for Intel Briefs and Reports. This enrichment system operates differently because it is anchored at the SPACE level. Rather than enriching entities independently, entity enrichment inherits context from the entity's primary space, internally described as "the chess piece must be on a board to enrich." In practice, this means a Workspace Intel Report may enrich against external context, and entity reports inside that space inherit that enrichment indirectly. The enrichment system is gated by scope-level controls and exposes status reporting back to the frontend so failed enrichment attempts can be surfaced gracefully to the user.

Architecturally, both systems are intentionally non-critical dependencies. Tavily enhances cognition quality when relevant, but the substrate, orchestration pipeline, and continuity system operate independently of it. If enrichment is unavailable, the system still functions normally using only durable cognition from the KRL substrate. The philosophy is deliberate: external RAG enrichment sharpens context, but it never replaces the structured continuity layer that forms the core of the system.

Layer [6]: Cognition Orchestrator

The orchestrator is the runtime brain of the system, but it does not generate the response itself. Its job is to assemble cognition: resolve policies, retrieve the correct continuity from the substrate, determine cognition depth, construct the structured prompt directives and citation rules, and hand that directive bundle to the route, which then streams the composer and language model response. This separation is intentional because it keeps the orchestrator testable and reusable, the same orchestration output can drive a text reply, voice interaction, or debugging replay without changing the cognition logic itself.

For each turn, the orchestrator follows a consistent sequence. First, it resolves the active cognition policy through a layered override system that determines things like retrieval depth, citation style, and consistency mode. Next, it retrieves the scope-correct continuity from the substrate and limits how much context is pulled based on the policy's retrieval depth. On cold-start sessions, it can also include temporary synthetic cognition as explicitly low-authority context without mixing it into durable memory. The orchestrator then classifies the user's intent into one of three cognition levels, Fast, Continuity, or Deep, using lightweight heuristics first and only using an LLM tiebreaker for genuinely ambiguous cases. Finally, it builds the cognition directives and citation rules that the composer will later use when generating the response.

This is where selective intelligence is enforced. The orchestrator's most important decisions are often the moments where it chooses to do less. A low-signal message like "I'm tired today" may trigger only a lightweight conversational response instead of deep pattern analysis and extensive retrieval. Even though the substrate and pattern engine could surface far more information, the orchestrator intentionally restrains the system from over-analyzing every turn. That discipline is what keeps the product feeling calm, relevant, and human instead of overly clinical or performative.

The orchestrator also manages the system's consistency behavior. After the asynchronous extractors and pattern engine finish updating the substrate, the orchestrator can refresh the continuity state and emit update markers so the UI silently syncs with the newly settled cognition state. In practice, this means the orchestrator is responsible not only for assembling cognition before the response, but also for closing the loop after the post-turn cognition pipeline completes.

Layer [7]: Composer

The composer is the final cognition layer before the language model. Its responsibility is not to decide *what* the system knows, that work has already been performed by the extractors, the KRL substrate, the pattern engine, and the orchestrator, but to decide how that cognition is translated into a coherent prompt for the model. The composer constructs the final LLM prompt by assembling the orchestrator's directive bundle, layering in the retrieved continuity, applying citation and evidence rules, attaching the appropriate safety overlays, and selecting the correct prompt skeleton for the active cognition scope and surface identity.

The composer operates as a rendering layer rather than a reasoning layer. It does not determine which patterns exist, which records are relevant, or how confidence was calculated. It only receives the structured cognition state that upstream systems have already resolved. This separation is deliberate. The composer knows that a TrustRebuilding pattern exists, but not how the pattern engine derived it. It knows that a record has high confidence, but not how retrieval scored it. By isolating prompt construction from cognition generation, the system keeps each layer independently testable, replaceable, and auditable.

Prompt construction itself is highly structured. The composer assembles multiple prompt regions into a single scoped instruction set: system directives, cognition directives, retrieved continuity blocks, pattern citations, safety overlays, scope identity instructions, tone constraints, and response policies. Different surfaces and cognition scopes receive different prompt skeletons. A USER-scope growth journal prompt emphasizes introspection and self-reflection. An ENTITY-scope prompt emphasizes relational dynamics and uncertainty discipline. A SPACE-scope prompt emphasizes systemic and environmental reasoning across multiple entities. The composer also controls citation discipline, whether the model should surface explicit evidence, balanced references, or implicit continuity depending on the active cognition policy.

Safety and behavioral governance are also enforced at the composer layer. Shared overlays, including crisis-response modules, anti-diagnosis rules, uncertainty framing, and contradiction-preservation directives, are injected consistently across all conversational surfaces. This ensures that cognition behavior remains stable regardless of which scope or interface the user is interacting with. The composer is therefore not merely “prompt engineering”; it is the policy-enforcement and cognition-rendering boundary between the structured substrate and the probabilistic language model.

WHAT THE COMPOSER SENDS TO THE LLM — AND HOW THE LLM UNPACKS IT

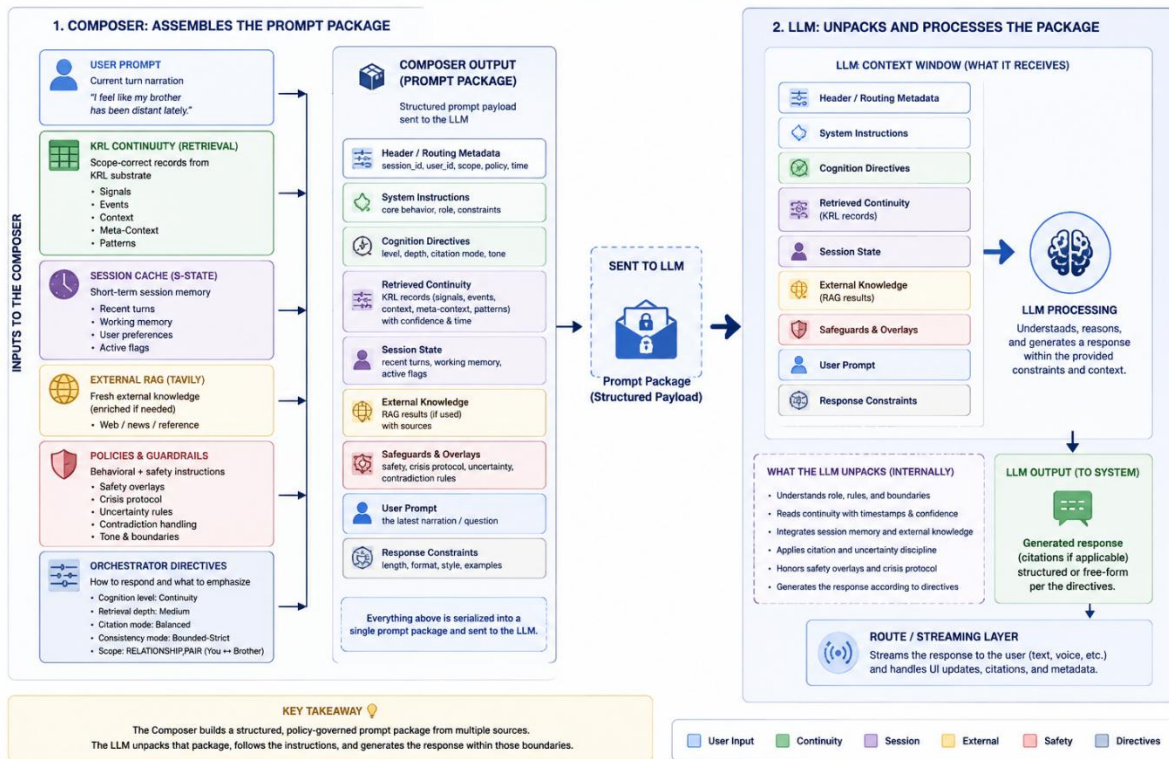


Figure V, Shows the process of the composer as it collects the inputs to construct the user prompts, retrieving context, rags, session memory, policies and guardrails and the prompt itself as it passes it on to the LLM.

Because the composer is isolated from the underlying cognition machinery, the language model itself becomes replaceable infrastructure rather than the center of the system. The orchestrator can assemble the same cognition bundle while different composers target different providers, model families, response styles, or modalities. A text stream, voice interaction, analytical report, or future multimodal interface can all consume the same substrate and orchestration output through different composition layers without rewriting the underlying cognition system. That separation of concerns is what allows the substrate, rather than the model provider, to remain the durable center of the architecture.

PART IV

Integrated Technical Additions: Substrate Semantics and Write Discipline

This section folds the technical-whitepaper depth into the unified paper. The architectural thesis is that persistent cognition is not extended chat history. It is a structurally separate substrate, written by extraction rather than generation, retrieved by deterministic addressing rather than embedding similarity, and governed by its own physics of decay, contradiction, and propagation. The LLM rents access to the substrate at the rendering edge; the substrate is the cognition.

The KRL as addressable cognition, not transcript memory

The Knowledge Representation Layer is organized as a 4 x 5 cognitive matrix. Four subject scopes define what a record is about: USER, ENTITY, SPACE, and RELATIONSHIP_PAIR. Five layer kinds define what kind of record it is: SIGNAL, EVENT, META-CONTEXT, CONTEXT, and PATTERN. This creates addressable cells instead of one large conversational blob. The orchestrator can retrieve "recent negative relational signals about this dyad" or "space-level meta-context drift" directly, rather than hoping similarity search returns the right memory.

Each durable record carries shared metadata that makes cross-layer compounding possible: layer, subcontext, dimension, subject_kind, subject_id, subject_aspect, perspective, polarity, temporal_scope, confidence, provenance, and evidence density. The important design choice is uniformity. A signal about a user, an event about a space, and a meta-context record about a relationship pair can be compared, filtered, scored, and traced using the same governance language.

Narration modes

A single user message can represent an observation, a narration, or a dyadic interaction. The extractor must distinguish these modes before it writes. If the user observes two people arguing, the system writes observed ENTITY records. If the user narrates a person from that person's pane, the system writes ENTITY records with the right perspective. If the user is a participant, the system may write RELATIONSHIP_PAIR records because the dyad itself is now the subject.

The extraction pipeline as the only durable write path

Extraction is the only durable write mechanism. The composer, analytics surfaces, and model output do not write durable cognition. A journal turn first produces the user-visible response, then the cognition tail runs extractors against raw user narration, persists proposals through the canonical upsert chokepoint, evaluates the pattern engine, bumps the cognition settle version, and records the trace. The product learns from the human, not from its own generated prose.

signal_extractor writes append-only SIGNAL records with polarity, behavioral markers, confidence, and decay semantics.

event_extractor writes EVENT records with participants, anchors, and links to earlier events.

context_extractor writes stable CONTEXT through contradiction and confidence gates.

meta_context_extractor writes drift-tracked interpretive state while preserving competing claims.

relationship_pair_extractor writes dyadic records keyed to the deterministic user-entity relationship identity.

space_filter_extractor decides which space context applies to the turn so routing is explicit rather than implicit.

Name-to-UUID discipline

Extractor prompts should let the LLM emit names and natural references, not invented UUIDs. Resolution to database identifiers happens in the persistence layer, where the system can verify that an entity exists, belongs to the current space when required, and maps to an allowed subject. This prevents a subtle but dangerous class of bugs where an LLM fabricates identifiers that look valid but point nowhere.

Clarification surfacing and contradiction handling

For stable facts and meta-context, the system treats contradiction as information rather than failure. A low-confidence contradiction becomes a clarification candidate; a high-confidence contradiction may update the active value while preserving history; a user-overridden fact blocks AI replacement until the user revisits it. The goal is not to make the substrate perfectly certain. The goal is to keep uncertainty visible and governed.

Integrated Pattern Engine Additions: Evidence, Lifecycle, and Safety

Four-tier escalation

Patterns should not surface because one sentence sounded meaningful. The pattern engine escalates from atomic records to recurring dynamics only when there is enough evidence density, confidence, temporal spread, lineage, and scope alignment. This is what separates signal intelligence from horoscope-style interpretation: the engine waits for accumulation.

Tier 1 patterns summarize local repeated signals inside one subject and layer.

Tier 2 patterns combine signals and events into a recurring dynamic inside a subject scope.

Tier 3 patterns connect multiple scopes, such as USER impact plus ENTITY behavior plus RELATIONSHIP_PAIR drift.

Tier 4 compound arcs connect patterns across time, spaces, and relationship pairs.

Pattern lifecycle: persisted state versus surface labels

The durable lifecycle should be understood as storage discipline plus computed surface language. The persisted lifecycle can remain simple, such as active or dormant, while the product surfaces richer labels like changing, repeating, weakening, reactivated, or expired based on evidence count, recency, suppression, reactivation history, and decay. This keeps the database state stable while allowing the UI to communicate what the user actually needs to know.

Evidence formula and lineage

Every pattern should be explainable through lineage. A compound pattern is not merely a name; it is the result of parent patterns and extracted records joining across dimensions. The Evidence

Chain walks from the surfaced pattern to `parent_pattern_ids`, from those parents to contributing records, from records to source messages, and from source messages to the actual user-authored language. Trust comes from traceability.

Pattern humanizer and safety charter

The humanizer converts internal detector names into user-readable language, but it must not upgrade a detector into a diagnosis. The safety charter is: do not diagnose the person, name the interaction pattern, show the evidence, use confidence, track contradictions, let the user correct the model, and do not fabricate if the entity has sparse data. This is both a product principle and a legal/ethical constraint.

Hallucination blocklist as a write-boundary backstop

A persistent cognition product must prevent invented pattern names from entering the durable bank. The runtime should filter weak or unsupported patterns before prompt injection, while the write boundary should reject known-fabricated pattern names and unsafe detector labels. The read-side gate prevents hallucinated cognition from being shown; the write-side gate prevents it from becoming substrate.

Cognition Architecture

The cognition architecture is a layered system that transforms user narration into structured continuity, extracting signals, events, patterns, and relationships into a governed substrate that allows the AI to reason with memory, context, and temporal continuity instead of replaying raw conversations.

Design Philosophy: Selective Intelligence

The single most important principle that has shaped our architecture is what we call selective intelligence. The instinct when building AI products is to do more on every turn: load more context, run more analyses, surface more patterns, generate longer responses. We have resisted that instinct deliberately.

Selective intelligence over maximal cognition.

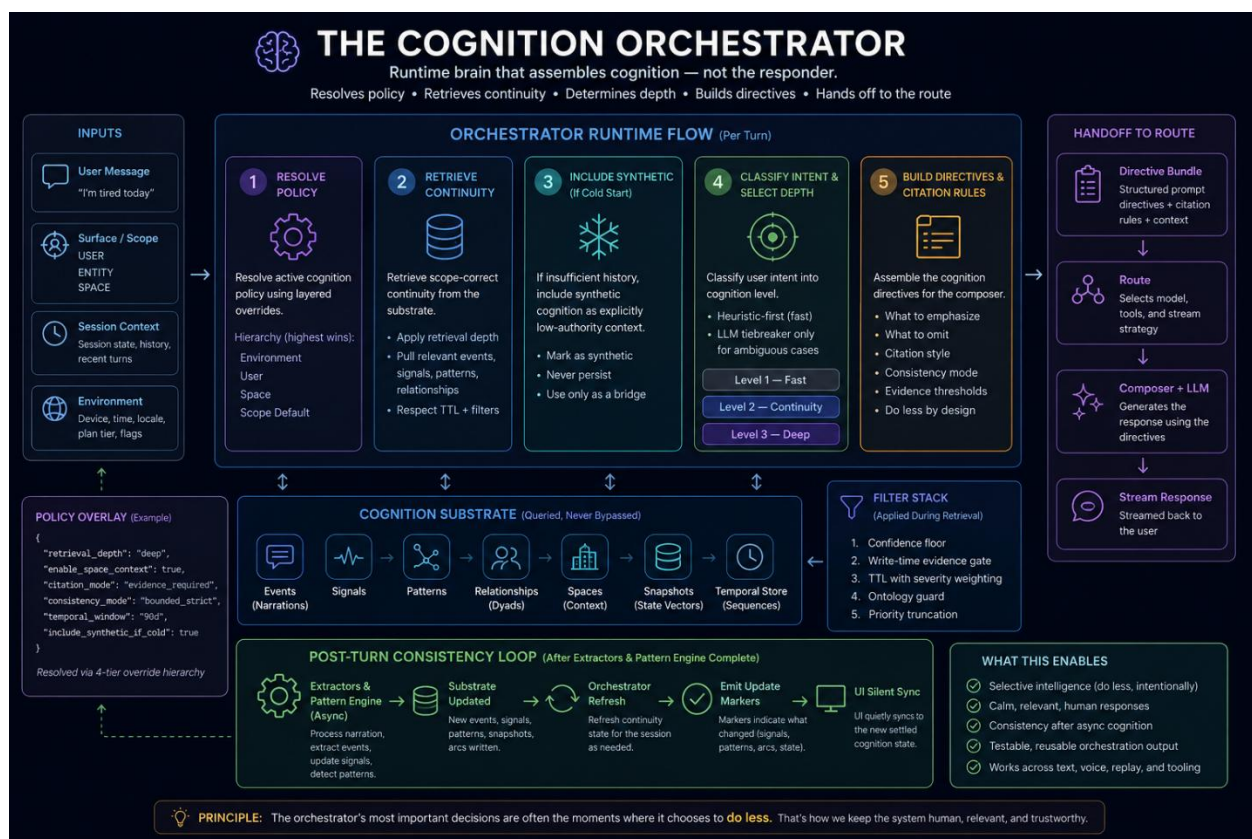


Figure W, Shows the cognition orchestration which steps from understanding the users prompt to determine the retrieval and effectively one needs to be engaged as we're profiling the prompt as it ultimately gets sent to the composer.

Selective intelligence means the system asks, before running anything expensive: does this query actually need it? When a user types "I am tired today," the right response is conversational and brief. Running pattern detection, citing five active dynamics, and surfacing two years of relationship history would be condescending, expensive, and bad product design.

When a user asks “what patterns do you see with my brother,” the right response is analytical and dense. Skipping pattern citations would be a failure of the product's core value proposition.

The Ontology Is the Product's Spine

Four scopes in the substrate. Four homes in the product. One coherent map.

Before walking through the screens, it helps to see how the substrate's ontology shows up in the product. Every cognition scope in the database has a corresponding home in the UI. Every UI surface reads from a scope-correct slice of the substrate. The two were designed together, not bolted to each other after the fact.

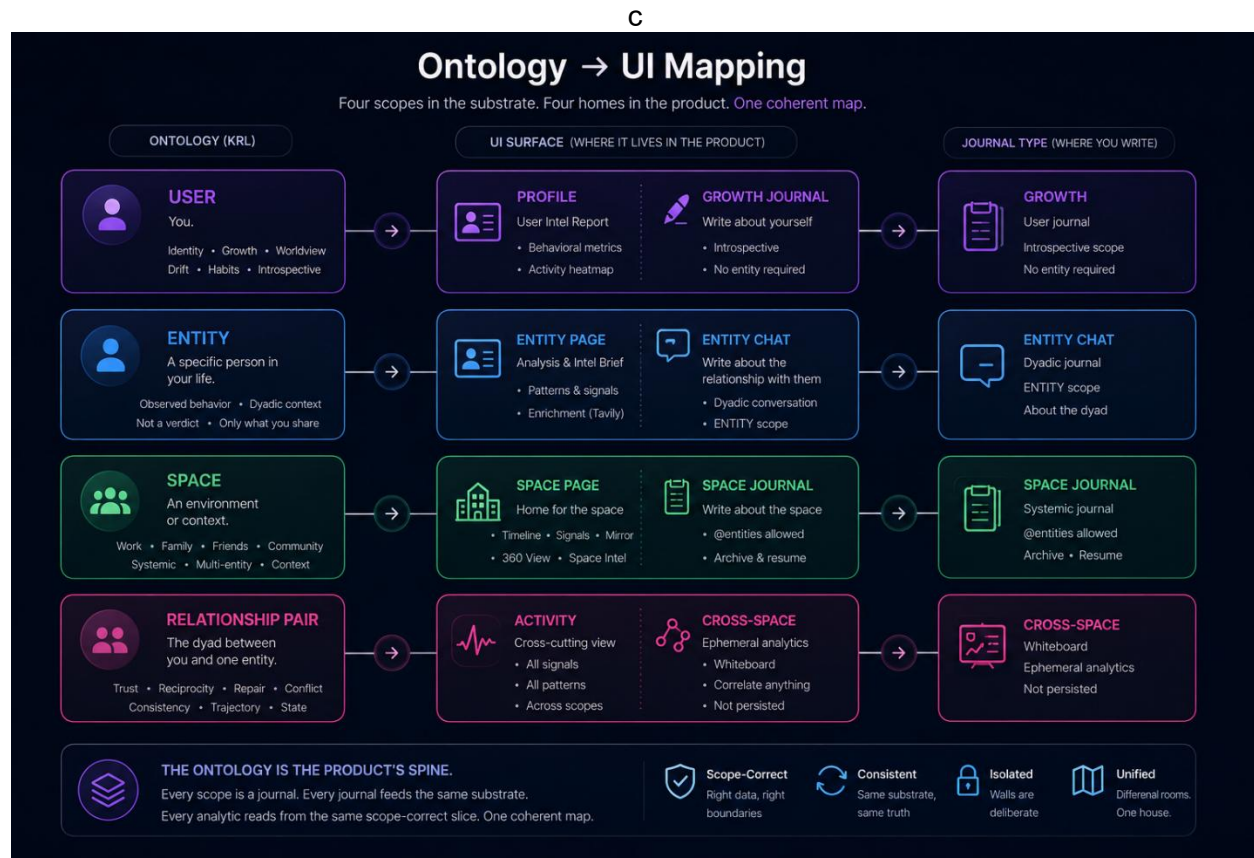


Figure X, Breaking down the domain specific journaling dimensions which is the pre-requisite for the KRL.

USER scope → Profile + Growth Journal

The USER scope holds the substrate's understanding of the person themselves, identity, growth, worldview, drift, habits. It surfaces in two places. The Profile is the read view: a User

Intel Report at the top with top signals and top patterns, an MBTI-style and attachment-style interpretation framed as Experimental, a behavioral metrics panel, and an activity heatmap. The Growth Journal is the write surface: a free-form introspective journal that requires no entity and no space. It is where you talk about yourself.

ENTITY scope → Entity Page + Entity Chat

The ENTITY scope holds observed behavioral material about a specific person in the user's life, never a verdict on the person themselves, only what has been narrated about them. The Entity Page is the analytical view: a pattern card per active dynamic, a signal trend chart, an Entity Intel Brief, and Tavily-backed enrichment where the system can add context from public sources. The Entity Chat (entity journal) is the dyadic conversation: writing about the relationship between you and them, in ENTITY scope, distinct from talking about them inside a space.

SPACE scope → Space Page + Space Journal + analytics

The SPACE scope holds the systemic view, a workplace, a family, a friend group, a community. It is the most complex scope because the entity overlaps are real and the signal space is the widest. The Space Page is the home for everything related to that environment: the active Space Journal, the Timeline of past sessions, a Signals view, the Mirror (where the AI asks the user clarifying questions to learn), the 360 view, and the Space Intel Report. Entities live inside spaces too, once assigned, they are referenceable from the space journal via @mention.

RELATIONSHIP_PAIR scope → Activity + Cross-Space

The RELATIONSHIP_PAIR scope holds the dyad-level state between the user and one entity, trust, reciprocity, repair, conflict, consistency, relational trajectory. It surfaces wherever the dyad is the right unit of analysis: most prominently in Activity (a cross-cutting view of all signals and patterns across scopes), and in the Cross-Space whiteboard (an ephemeral analytic surface for running ad-hoc correlations across imported spaces and entities).

Three Cognition Scopes

USER, ENTITY, SPACE. The same substrate. Different identities. Different rooms in the same house.

Humans do not carry continuity in a single bucket. The way you think about yourself is different from the way you think about your sister, which is different from the way you think about your workplace. We modeled the substrate after that distinction. Three cognition scopes, each with its own identity, sharing the same infrastructure underneath.

USER scope - introspective

When the cognition is about you, it speaks in second person, it draws on identity, growth, habits, worldview, and internal drift. It does not pivot to talking about an external entity unless you bring one in. The framing is reflective, not analytical.

ENTITY scope - dyadic

When the cognition is about another person, it speaks about the relationship between you and them. Trust, reciprocity, repair, conflict, consistency, relational trajectory. The subject is the dyad, not the entity in isolation, because we have no way of knowing the entity outside what you have told us. We are careful to never produce verdicts about another person. We confirm patterns we have observed when the evidence supports them. We never confirm conclusions.

SPACE scope - systemic

When the cognition is about an environment, a workplace, a family, a friend group, a community, it speaks about group dynamics, organizational stress, social structure, environmental drift, multi-entity interactions. This is the most complex of the three because the signal space is wider and the entity overlaps are real. It is also the most distinctive, there is no consumer product we have seen that does this level of space-scoped cognition.

Three rooms in the same house. The walls are deliberate. We will never collapse them into a single feed, because the room you are in changes what continuity is appropriate.

Cognition Scopes versus Surfaces

A common point of confusion when explaining the system is the difference between a cognition scope and a UI surface. The two are related but distinct, and the architecture depends on keeping them separate.

Three cognition scopes

There are three cognition scopes, each with its own scope identity, prompt skeleton, retrieval shape, and policy defaults: USER (introspective), ENTITY (dyadic), and SPACE (systemic). Every journal turn happens in exactly one scope. The orchestrator resolves the scope at the start of the turn and the rest of the pipeline obeys it.

Ten or so UI surfaces

The product has roughly ten places where the user can interact with cognition: the space journal, the user journal, the entity chat, the growth journal, the Mirror surface, the 360 view, the cross-space whiteboard, the onboarding context conversation, the analytical questions over the substrate, and the burst surface for over-cap usage. Some surfaces are journals (write paths into the substrate). Some are analytics (read-only views). Each surface declares which scope it is in, and the orchestrator inherits the scope identity from that declaration.

Why this matters

A scope is an identity. A surface is a place. The same scope can have multiple surfaces (USER scope has both the Growth Journal write path and the Profile read path). A single surface always lives in one scope (the Entity Chat is always ENTITY scope, never USER). Keeping the two concepts separate is what lets us add new surfaces without changing how cognition is composed, the surface just declares its scope, and the orchestrator handles the rest.

Runtime Cognition Governance

Cognition Policies

A four-tier override hierarchy that lets every surface tune itself without changing code.

Different surfaces in the product need different cognition apertures. The space journal wants more context. The entity journal wants the dyad. The user journal wants introspective depth. The growth journal wants narrative continuity. We did not want to fork the orchestrator per surface. So, we built a cognition policy layer.

A policy is a small JSON contract: `retrieval_depth`, `enable_space_context`, `citation_mode`, `consistency_mode`, and a few other flags. Policies resolve through a four-tier hierarchy: the scope's default policy, the space-level override, the user-level override, and the environment-level override. Environment beats user; user beats space; space beats scope default. The orchestrator resolves the effective policy at the start of every turn.

This lets us A/B-test cognition behavior per cohort without redeploying code. It lets us tune a specific user's policy for a debugging conversation. It lets us shape how aggressive the system is at retrieving context for a high-stakes space. The four-tier hierarchy turns out to be the right number, three felt brittle, five felt over-engineered.

Cognition hierarchy

Operationally, selective intelligence is implemented as a three-level cognition hierarchy.

Level [1]: FAST PASS

The default for ordinary low-signal narration. Lightweight summary, basic signal extraction, minimal continuity lookup. No deep pattern lineage, no predictive trajectory, no relationship-state evolution. Responses are short and respect the user's energy.

Level [2]: CONTINUITY PASS

Fired when the user mentions a specific entity, narrates a recent event with emotional weight, or engages with a recurring topic. The system retrieves related events, unresolved arcs, dyadic continuity, and active pattern outputs. Responses honor history without becoming clinical.

Level [3]: COGNITIVE PASS

Fired only when the user explicitly invites analysis, asks for instances, requests pattern lineage, or surfaces something that warrants major retrieval. The full pattern engine runs, evidence is cited with confidence framing, and the response engages the user's analytical intent.

Classification is Heuristic First

a small regex layer that runs in roughly a millisecond, with an optional LLM tiebreaker only for ambiguous mid-band messages. Heuristics handle the vast majority of cases. The tiebreaker is reserved for genuine edge cases where the user's intent is genuinely uncertain.

The cost of ignoring this principle

Early in development we ran every turn at maximum cognition depth. The result was a system that produced thoughtful, dense analyses for users who had just said “feeling alright today.” The responses were technically impressive and experientially wrong. They read as if the AI was straining to prove it was paying attention. Users reported it felt performative.

After implementing the cognition hierarchy, the same low-signal messages began receiving short conversational responses, and the analytical depth was reserved for moments when users were actually asking for it. The system felt calmer and more present. This is the kind of outcome we now optimize for: not the maximum amount of intelligence visible per turn, but the right amount given what the user actually wants.

Implications throughout the stack

Selective intelligence is not a single feature, it is a discipline that shapes every layer. The retrieval layer asks whether continuity is needed before fetching. The pattern engine asks whether a pattern is worth surfacing before injecting it. The composer asks whether a citation is warranted before constructing it. The cognition policies, which we describe later, are the mechanism by which different surfaces declare their cognition aperture. Every dial we have built is in service of letting the system do less, more accurately, more often.

Async Cognition Consistency

Eventual / bounded-strict / strict. The pipeline acknowledges that cognition takes time.

The user's LLM response streams back fast, because that is what the user is waiting for. The cognition pipeline that processes the user's narration, extractors, pattern recomputation, snapshot updates, runs after the stream finishes. If the user re-asks immediately, the new substrate state has not settled yet. The model would answer based on stale patterns. That is a bad UX moment.

We modeled three consistency modes per cognition policy. Eventual is the default, the response streams immediately and the substrate updates in the background. Bounded-strict emits a settle-version marker over the same SSE stream when the post-turn cognition has finished, and the UI silently refreshes. Strict blocks the response until the substrate has settled, slower but guaranteed-fresh, useful for analytical questions where the user is explicitly asking for the latest state.

Bounded-strict is the most interesting of the three. It removed the re-ask cliff from the product without forcing every turn to wait for the slowest layer.

Synthetic Cognition - the Cold-Start Layer

How we handle the first day, before the substrate has had time to learn.

A new user has no substrate. The pattern engine has nothing to fire. The orchestrator has nothing to retrieve. The cold-start problem is real, and the failure mode of most memory systems is that they over-promise on day one and then under-deliver as the model invents continuity it does not have.

We solved this with a layer we call synthetic cognition. During onboarding and the first few sessions, a registered synthesizer writes synthetic records into a separate session-scoped blob, not into the durable substrate. Those synthetic records are tagged with LOW authority. They participate in retrieval, but the model is explicitly told that they are bootstrap scaffolding, not observed continuity. As the durable substrate fills out, the synthetic blob's relative weight drops. We call this the hourglass: synthetic cognition tapers as durable cognition accumulates.

The discipline that matters here is that synthetic cognition can never become durable cognition. The two are stored separately, retrieved with different weights, and tagged differently in the prompt. We do this because a memory layer that recursively ingests its own outputs eventually drifts into hallucination, a kind of epistemic feedback loop. The asymmetric authority is the firewall.

Compounding Filter Stack

Five filters between raw narration and the prompt. Each strips a different kind of noise.

Signal-to-noise is the whole game. We built five filters that sit between what a user types and what reaches the LLM. None of them is novel on its own. The discipline is in stacking them deliberately, so each one removes a class of noise the next one cannot.

The compounding part is the point. Filter five only has to do a little work because filters one through four already did most of it. None of them is impressive alone.



Figure Y, Five filters stacked between narration and prompt. Confidence floor, write-time evidence gate, severity-weighted TTL, read-time ontology guard, priority truncation. Each strips a different class of noise.

Filter [1]: Confidence floor

When an extractor produces a signal, it also produces a confidence score. If the score is below a configurable floor, the signal never enters the substrate. It sounds blunt, and it is, but it removes the largest single class of noise: low-confidence LLM extractions that are technically structured but semantically thin.

Filter [2]: Write-time evidence gate

A pattern only persists when enough signals across a real time window agree. We do not fire a pattern off a single utterance. This is the difference between a system that detects mood and a system that detects dynamics. It also stops the second-most-common class of noise: a single emotionally charged narration triggering a spurious pattern that would then influence future turns.

Filter [3]: TTL with severity weighting

Older signals decay faster than recent ones. Critical signals decay slower than minor ones. We treat decay as severity-weighted, not flat. A signal of withdrawal from three weeks ago carries less weight than one from yesterday, unless yesterday's signal is part of a longer pattern, in

which case the older one matters again. We also do not delete dormant signals; we hide them. If they become relevant again, the system can pull them back.

Filter [4]: Read-time ontology guard

When the orchestrator retrieves patterns for the current turn, it can only retrieve from a closed vocabulary of patterns that exist in the bank. If the system has no patterns for this user's scope, it does not hallucinate categories to fill the silence. It tells the model so. This is the rule that stops the most pernicious failure mode of LLM-driven memory systems: inventing patterns that sound plausible because the prompt asked for them.

Filter [5]: Priority truncation

Even after the first four filters, the orchestrator still trims the retrieved context to a small budget before composing the prompt. Critical patterns first. Recent signals next. Everything else cut. The point is that the prompt that reaches the LLM is small, dense, scope-correct, and cited. The LLM is doing rendering, not reasoning over a hundred turns of history.

Temporal Continuity Cognition Layer

the temporal continuity layer that adds intra-turn sequencing, narrated anchors, and arc detection to the existing pattern engine.

What temporal continuity adds

The temporal continuity layer adds arc detectors that fire when a recognizable sequence of patterns occurs over time, and that surface the trajectory rather than just the latest state.

Intra-turn sequencing

A single narration often contains multiple temporally-ordered events. “This morning in standup, Mike interrupted me. Then at lunch, Sarah brought it up unprompted. By end of day, Mike actually apologized.” The event extractor now emits `temporal_order` and `narrated_anchor` fields so the substrate captures the sequence, not just the events. Downstream pattern detectors can then reason about “what came before what.”

Anchor parsing

Users narrate in natural temporal language, “last Tuesday,” “that first weekend,” “a month after the wedding.” The `narrated_anchor` parser resolves these to wall-clock estimates with confidence, so the substrate can answer questions like “what was happening around the time of the wedding” without forcing the user to enter calendar dates.

Arc detectors

Four arc detectors are shipped and live: RepairThenRelapse, TrustRebuilding, EscalationArc, and WithdrawalSpiral. Each reads the substrate's temporal sequence of relevant signals and fires when a recognizable trajectory has occurred. Arcs surface as a different kind of pattern, narrative rather than aggregate, and the composer is instructed to honor the trajectory framing when surfacing them.

Event correlation index and relationship state snapshots

Two background services support the temporal layer. The event correlation index resolves implicit references across narrations, when a user says “the thing last week,” the resolver attaches the right earlier event id so the substrate knows what “the thing” refers to. Relationship state snapshots write periodic state vectors for each user-entity dyad, trust, reciprocity, repair, conflict, consistency, relational trajectory, so that the substrate can answer “how was this relationship doing at this point in time” without re-deriving from raw signals.

Why this is the premium-tier wedge

Static pattern detection is the baseline cognition product. Temporal continuity is the cognition product that compounds in value the longer a user is on the platform. Three months of substrate produces arcs. Six months produces inflection points. Twelve months produces longitudinal trajectory. The longer the user stays, the more uniquely valuable the substrate becomes. We expect temporal continuity to be the feature that converts free-tier curiosity into paid-tier commitment.

Observability, Versioning, and Experimentation

Once we had three cognition scopes, three policies, and an orchestrator that could resolve them, we needed a way to evolve any of those things safely. The versioning and experimentation layer is the deployment-safety substrate that makes cognition iteration possible.

Per-layer version constants

Every cognition layer carries a version constant: orchestrator, pattern engine, signal extractor, composer. We record all four on every turn's trace. When a behavior changes ships, the version constant bumps. Trace consumers can then filter by version when analyzing outcomes, we can ask “what is the citation rate when pattern engine is at v2.0 versus v2.1?” and get a clean answer. Without per-layer versioning, behavior regressions would be invisible until users complained.

Deterministic variant assignment

We assign users to experiment variants deterministically: a hash of (user identifier, experiment name) maps to a bucket between zero and ninety-nine, and the bucket determines which variant the user sees. The same user always lands in the same variant for the same experiment. Changing the rollout percentage shifts the population boundary, but every user's membership remains stable up to that boundary. This avoids the most common A/B pitfall, where users flip between variants and the experiment data becomes unusable.

Variant policies

Experiment variants resolve to Cognition Policy instances. This means an experiment is really “some users get policy A, others get policy B,” and the orchestrator already knows how to consume different policies. We do not have to thread experiment context through every layer of the system; the experiment's output is a variant policy, and the orchestrator uses it the same way it uses any other policy.

What this unlocks

This infrastructure is the prerequisite for any post-launch cognition iteration. Without it, every behavior change would be a monolithic production cutover with no rollback path. With it, we can introduce a candidate policy, expose it to a small fraction of users, measure outcomes against the control group through trace queries, and promote or roll back based on data rather than intuition. The infrastructure is in place even though the experiment registry is currently empty, we deliberately deferred the first experiment until real users are in the system, because experiment design without real users is mostly theater.

PART V

Operational Insights

What instrumenting every cognition turn taught us, including the moments when the system surprised us, persistent memory can coexist with session memory.

The Trace Layer (and what it taught us)

Every cognition run writes a structured trace. We have learned more from these traces than from anything else.

Each turn writes a structured trace describing what the orchestrator decided, what the retrieval layer fetched and why, what patterns the engine surfaced, what the composer assembled, and what the LLM said. The trace is not for the user, it is for us. It lets us debug behavior weeks after the fact without having to re-ask the model what it did.

Several of the most important architectural decisions in this paper came directly from things the trace layer caught. We will name a few because we owe the reader honesty.

Hallucinated pattern names

Early in development we saw the LLM citing patterns that did not exist in our bank. The model was inventing plausible-sounding category names because the prompt left room for it. We added the read-time ontology guard (Filter 4) and a sparse-substrate directive that tells the model to acknowledge silence rather than fill it. The traces confirmed the hallucination stopped.

Stream hang from pattern enrichment

A subtle bug had every pattern in a turn triggering a per-pattern LLM call to generate emotional alternatives, serially, with no outer timeout. Long turns would hang. The trace layer showed the wait clustered around the enrichment step. We wrapped the call-in a per-call timeout and moved an inexpensive type-check pre-LLM. Latency dropped immediately.

Settle versions and the re-ask cliff

We watched users re-ask questions immediately after the AI responded, because they could tell the response was based on stale pattern state. The pipeline was running fine, extractors and pattern recomputation just happened after the LLM streamed back. We added a settle-version marker that the stream emits when the post-turn cognition has finished. The UI now waits for the marker and refreshes silently. The re-ask cliff went away.

We have been wrong in instructive ways. Cataloguing those mistakes honestly is part of what builds a substrate worth trusting.

System Flow: Extract versus Retrieve

A common question about systems like ours is: how does the system know whether to add to its understanding or query it? The answer involves a deliberate separation between asynchronous extraction (which adds) and synchronous retrieval (which queries). This separation lives along the timeline of a single turn.

The turn lifecycle

Synchronous path, visible to the user

- The user types a message into the journal.
- Authentication and quota enforcement run in milliseconds.
- The orchestrator runs intent classification (~1ms heuristic, optional LLM tiebreaker only on ambiguous L2 default cases).
- Based on classification, the orchestrator decides what continuity to load and what pattern injection to assemble.
- The composer produces a prompt and streams the LLM response to the user.
- The user-visible response completes. Time elapsed: typically, two to five seconds, mostly LLM generation.

Asynchronous path, invisible to the user

- After the user message persists, multiple extractors fire in parallel.
- Signal extractor reads the message for behavioral signals and writes them to KRL.
- Event extractor reads the message for discrete narrated events and writes them to KRL.
- Meta-context extractor reads for drift-prone interpretive reads and writes them to KRL.
- Relationship-pair extractor specifically looks for dyad-level dynamics.
- The pattern engine receives the new signals, re-evaluates relevant pattern templates, and writes updated pattern records.
- Embedding generation indexes the message for future similarity queries.
- These steps continue running after the user has received their response. They do not block the user-visible turn.

How the system knows which mode

The extraction path always runs on every user message, we want every narration to feed the substrate. The retrieval path runs as part of the synchronous turn and is selective. The orchestrator decides what to retrieve based on the intent classification, the user's explicit @-mentions, and the heuristic signals in the message itself. A query like “what patterns do you see with my brother” triggers an analytical retrieval that surfaces patterns and signals; a narration

like “my brother apologized for the thing last week” triggers normal-depth retrieval and is treated primarily as input to extract from on the asynchronous side.

Why we separated them

Keeping retrieval synchronous and extraction asynchronous serves both performance and product semantics. The user does not wait for the pattern engine to recompute before seeing a response. The pattern engine does not have to produce results in user-perceptible time. And the substrate accumulates monotonically across every turn, regardless of whether retrieval surfaced anything. The result is a system that feels fast on the read side and gets smarter on the write side, every turn, without the user noticing the work.

PART VI

Integrated Runtime Additions: Orchestration, Policy, Continuity, and Synthetic Cognition

Orchestrator contract

The orchestrator is the runtime brain that resolves scope, policy, cognition depth, continuity retrieval, prompt assembly, streaming, trace capture, and post-turn consistency. It should return a structured result, not just prose: response text, token usage, referenced substrate, cognition level, trace metadata, terminal frame payloads, and any pending background cognition status.

Cognition routes and levels

The product uses different cognition depths as both a UX discipline and a unit-economics discipline. Fast Pass handles low-signal turns without expensive retrieval. Continuity Pass loads the relevant substrate slice for ordinary context-aware dialogue. Cognitive Pass performs deeper pattern retrieval, evidence review, and cross-space reasoning when the user explicitly asks for analysis. The classifier should be heuristic-first, with an LLM tiebreaker only for ambiguous cases.

Policy governance

CognitionPolicy should be treated as a runtime control plane, not a small settings object. It governs history depth, confidence floors, anomaly thresholds, consistency mode, activation of meta-context and space dynamics, analytic voice, correction handling, relationship propagation,

background cognition, and surface-specific retrieval behavior. The value is not that every field is user-visible; the value is that the system can change behavior without rewriting the substrate.

Continuity block and token budget

The continuity block is the read-side assembly layer. It pulls durable records, active patterns, synthetic cold-start scaffolds when permitted, recent messages, and surface-specific instructions into a bounded prompt payload. Token budgeting is part of cognition quality: the goal is not to send everything, but to send the fewest high-authority records that preserve meaning. Priority truncation should drop low-authority and low-relevance material before it drops durable, recent, high-confidence evidence.

Synthetic cognition with asymmetric authority

Synthetic cognition solves cold-start, but it must never become durable substrate. Its authority is asymmetric by design: durable evidence appears first, synthetic context appears later, the prompt labels synthetic material as tentative and lower-authority, and the retrieval cap keeps synthetic content small relative to durable context. Synthetic material may help the LLM ask better early questions; it cannot override user-authored durable evidence.

The four structural defenses against regenerative learning

Extractors receive raw user narration, not the LLM response.

Durable extracted_records require user_narration provenance at the persistence chokepoint.

Synthetic cognition lives in a conversation-scoped store with separate provenance and lifetime.

Runtime prompt assembly labels DURABLE and SYNTHETIC blocks separately so the LLM cannot treat them as equal authority.

The highest-level invariant is simple: the KRL must never recursively ingest itself.

Legal, Ethics, and the Crisis Protocol

The things we put in the product because we built it for adults living adult lives.

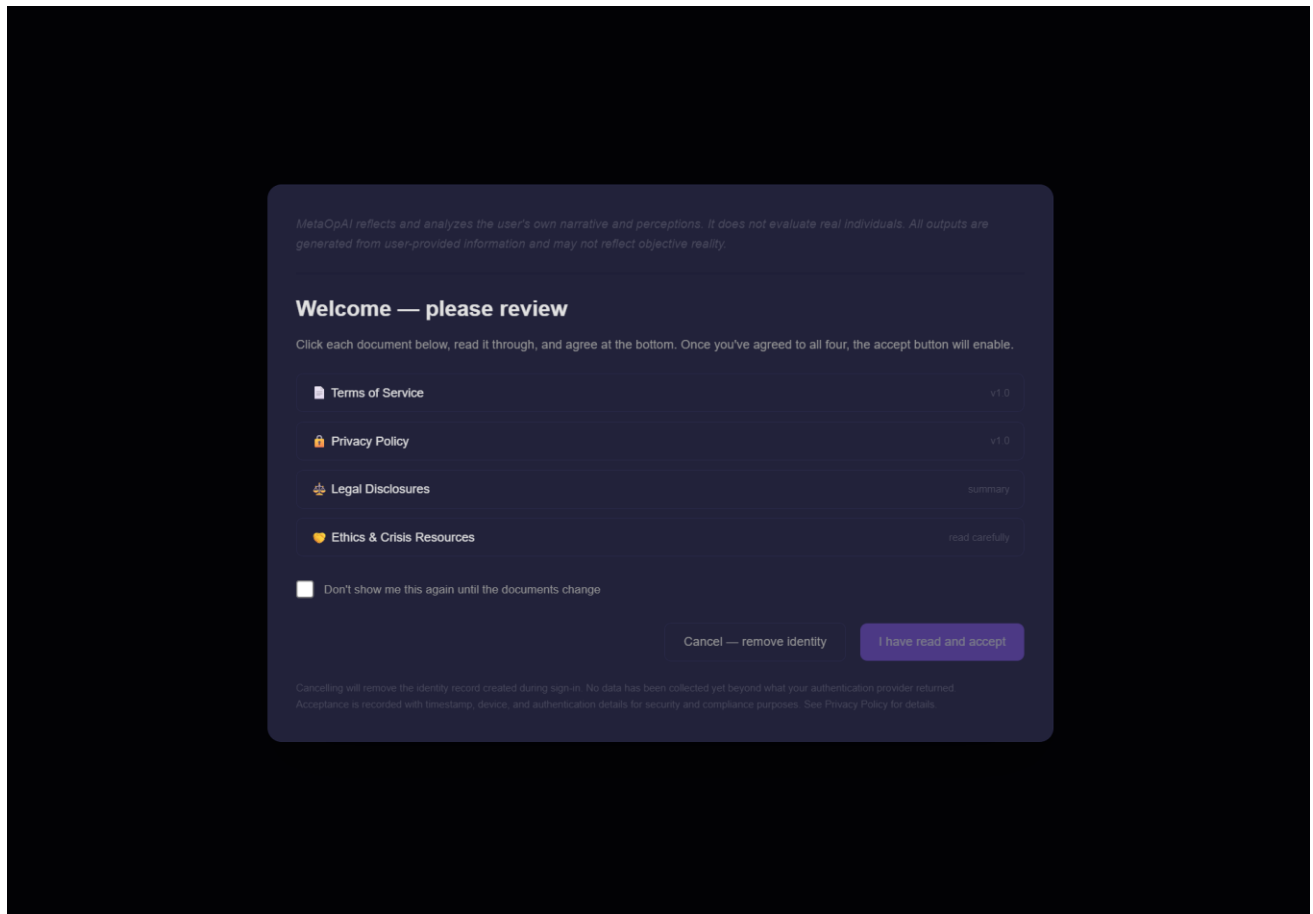


Figure Z, The legal, tos, private policy, and ethics legal hold before a user can even sign up and will prompt the user if a change was made requesting acceptance again which is logged from their ip address and local time.

What this product is, and is not

MetaOpAI is a signal intelligence product. It is not therapy. It is not medical or psychological advice. It is not a diagnostic tool. It is not a substitute for professional care, legal counsel, or close human relationships. We say this in plain language in the Terms of Service, in the Privacy Policy, in the Ethics Statement, and on every surface where a user might mistake an interpretation for an authoritative claim. The MBTI-style and attachment-style outputs are labeled Experimental for the same reason. They are pattern interpretations, not assessments.

How we handle high-stakes content

We do not produce verdicts on other people. The system can confirm a pattern it has observed when the evidence supports it. It will not confirm a conclusion. When a user is in a high-stakes relational context, suspicion of betrayal, marital crisis, a workplace harassment situation, the cognition policy shifts to counter-weighted reflection: the substrate retrieves balanced material,

the composer is instructed to surface multiple interpretations, and a Terms-of-Service clause makes explicit that nothing the system produces is admissible as evidence or a substitute for legal or professional judgment.

The crisis protocol

Suicide, self-harm, and violence signals override every other AI rule in the product. There is one shared safety module, imported by every conversational surface, that handles the routing. When a crisis signal is detected, the system surfaces the appropriate resources, 988 (Suicide & Crisis Lifeline), Text HOME to 741741 (Crisis Text Line), 911 for emergencies, and the composer is instructed to respond supportively without diagnosing. The crisis resources live in one place in the code so a resource change happens in one place across all ten chat surfaces. We do not bury them in a footer. They are not an afterthought.

Data privacy

User-owned data is our first cross-cutting principle. The user can export their entire substrate as JSON at any time. The user can delete their account at any time. Behavioral capture is gated by an explicit collection toggle. Federated authentication only, no passwords stored. Encryption in transit and at rest. The trace layer is for internal debugging and is not exposed to third parties. The Privacy Policy spells out exactly what we collect, why, and for how long. We try not to use the word minimal because the word does not mean anything; we tell the user what is stored and let them judge.

PART VII

Integrated Operational Additions: Security, QA, and Runtime Hardening

Azure-oriented deployment posture

The deployment model favors managed cloud primitives: Azure Container Apps or equivalent PaaS compute, Azure Database for PostgreSQL, Azure Cache for Redis, Key Vault, Managed Identity, centralized logging, private networking where appropriate, and provider abstraction for OpenAI / Azure OpenAI. The substrate must remain durable even if models, prompts, retrieval policies, or frontend surfaces change.

Authentication, consent, and privacy posture

The platform should continue treating authentication, consent logs, data export, deletion, and provider isolation as first-class system concerns. MetaOpAI handles highly personal narrative data. Trust depends on giving users clear control over what is captured, how it is used, and how it can be deleted or exported.

Cognition lock and consistency

The cognition lock exists to prevent overlapping write-side cognition from corrupting settle versions or racing pattern updates. In distributed mode, Redis-style SET-NX-PX semantics with holder-only release are the correct primitive: one worker acquires a lease, the lease expires safely if the worker crashes, and a release script prevents one worker from accidentally releasing another worker's lock.

QA harness and validation

The QA harness should validate more than whether an API returns 200. It should test sparse-entity behavior, contradiction gates, hallucinated pattern-name rejection, prompt payload construction, settle-version behavior, cross-space retrieval, Evidence Chain traceability, and high-stakes silent-skip rules. A cognition product needs QA against epistemic failure modes, not only functional failure modes.

Vision and Business

Why we built this, who it is for, how it makes money, why we think it lasts.

Vision: Be Understood

The product strategy is the same as the human need.

Our motto, Be Understood, is also our product strategy. We believe the deepest human need an AI can serve is the experience of being seen accurately over time. Not flattered. Not therapized. Not advised. Just seen, with the kind of patient continuity that a long-time friend brings, where the context is already there, the patterns are already noticed, and the questions can start somewhere later than the beginning.

Most consumer AI products today do not solve for being understood. They solve for being clever in the moment. The moment ends, the cleverness evaporates, and the user starts the next conversation from zero. We are building toward the opposite shape, and we say more about why that shape compounds in value over time in the Defensibility chapter.

This is not the same as memory in the trivial sense. Memory is necessary but not sufficient. The product also needs to recognize, to distinguish a one-off from a pattern, a steady-state from an inflection, a contradiction from a complication. Recognition is what turns memory into understanding. The substrate we are building is recognition-first; memory is the substrate underneath it.

If we get this right, MetaOpAI becomes the place a person goes to think about their actual life, the people in it, the spaces they move through, the patterns they cannot quite see from inside themselves. Not as a substitute for human relationships. As a way to be sharper inside them.

Operating Principle: Clarity

If Be Understood is the product strategy, Clarity is the operating principle. Every dial we have built is in service of clarity, for the user, for the system, and for us.

Clarity for the user

Cited evidence on every pattern. Confidence scores visible. Provenance traceable. The user can ask why the system surfaced something, and the answer is in the data, not in a black box. We do not produce conclusions; we surface patterns with evidence and let the user form their own conclusions. We do not hide what the system saw; we show it.

Clarity for the system

Every cognition turn writes a structured trace. Every layer carries a version. Every record carries provenance. Every threshold is defined in one place. When something behaves unexpectedly, the audit trail is recoverable. We do not debug by guessing what the model did; we debug by reading what the trace recorded.

Clarity for us

The architecture is layered, separable, and versioned. When we want to evolve the pattern engine, we change the pattern engine. When we want to retune the orchestrator, we retune the orchestrator. We do not have to reason about the entire system to change a part of it. That is the clarity we owe future-us, and the clarity that makes shipping the next thing possible.

Market and Competitive Positioning

The market we are entering

Consumer AI is in a wrapper-saturated phase. Most products sit between an LLM and a UI with minimal substrate underneath. The competition is on UX polish and prompt cleverness, not on cognition depth. We believe this phase is temporary. As model providers commoditize and users stay on platforms longer, the discriminator shifts from “which model does this product use” to “which product remembers me best.” That shift is the market opening we are building into.

Adjacent categories and why we are not them

Journaling apps (Day One, Stoic, Reflectly) do free-form text entry without structured cognition. We do structured cognition with the journal as the ingestion surface. Therapy apps (BetterHelp, Talkspace) connect users to licensed professionals. We are not a therapy product and we say so explicitly. AI assistants (ChatGPT, Claude, Pi, Replika) are conversational with shallow memory. We are signal intelligence with a substrate.

The competitor that worries us most is not any of the existing categories. It is the major model providers shipping their own cognition layer. OpenAI's persistent memory feature, Anthropic's Claude Projects, and Google's Gemini memory are all early moves in that direction. Our bet is that a substrate-first product, built with the discipline a general-purpose AI cannot afford to apply, will produce a better experience for the specific human-relational domain we focus on, and that focus is the defensible edge.

Why we focus on three subject domains

Self, others, and environments are the three places humans carry continuity in their actual lives. A cognition product that handles all three well is genuinely useful. A cognition product that handles all of them generically, try to track everything from your taxes to your fitness to your career to your kids' homework, is shallow on all of them. We chose depth over breadth. The domain shape (self, others, environments) is broad enough to be a real product and narrow enough to be done well.

Integrated Investor Additions: Market Position, Business Model, and Defensibility

Why AI memory architecture is the business thesis

The investor thesis is that LLMs are becoming more interchangeable while memory architecture becomes more valuable. Many AI products monetize repeated token replay: the user pays the model to rediscover context the product should have structured and stored. MetaOpAI moves value from transient prompt length into a persistent cognition substrate. That is both a product advantage and an economic advantage.

Market position

MetaOpAI is adjacent to journaling, therapy-adjacent reflection, personal CRM, knowledge management, quantified self, and general AI assistants, but it is not identical to any of them. The category claim is private signal intelligence for human life: a system that converts lived narration into structured, evidence-backed continuity across self, people, relationships, and environments.

Business model logic

The subscription model monetizes compounding continuity rather than one-off answers. A user who journals for one day receives reflection. A user who journals for months receives pattern memory, evidence trails, cross-space correlation, timeline reconstruction, and relationship-state movement. This supports tiering by spaces, entities, token budget, advanced surfaces, and burst capacity without turning the user experience into a raw per-token meter.

Defensibility

The moat is not a prompt. It is the combination of substrate design, write discipline, scope-specific ontology, pattern lineage, read/write separation, safety boundaries, and user-owned compounding data. Competitors can copy a UI faster than they can copy months of governed per-user substrate and the engineering discipline required to keep it clean.

Validation and roadmap

Validation should be framed across three tracks: technical validation, product validation, and business validation. Technical validation proves that extraction, patterning, retrieval, traces, and safety gates work. Product validation proves that users understand and repeatedly use the surfaces. Business validation proves that compounding continuity creates willingness to pay and retention. The roadmap should continue from private beta to broader beta, mobile distribution, improved social-media onboarding, voice-to-text, whiteboard-style cognition, and deeper analytics once usage data confirms the strongest surfaces.

Investor takeaway

MetaOpAI is not betting that one model will win. It is betting that users will increasingly need durable, private, structured cognition that survives model churn. If the LLM is the rented engine, the KRL is the owned chassis. The company becomes stronger as the substrate compounds.

Business Model

Subscription tiers

We charge a monthly subscription with three paid tiers and a Free tier. Tiers differ by token allowance, resource caps (spaces and entities), and burst purchasing. Feature gates are flat, the substrate is the product, and the substrate is for everyone. Pricing is unified across web, iOS, and Android at \$11.99 / \$22.99 / \$44.99 per month. Daily token caps reset at the user's local 6 a.m. so a heavy journaling week does not run out of budget. Resource caps graduate from one space and two entities on Free up to fifteen spaces and thirty entities on the top tier.

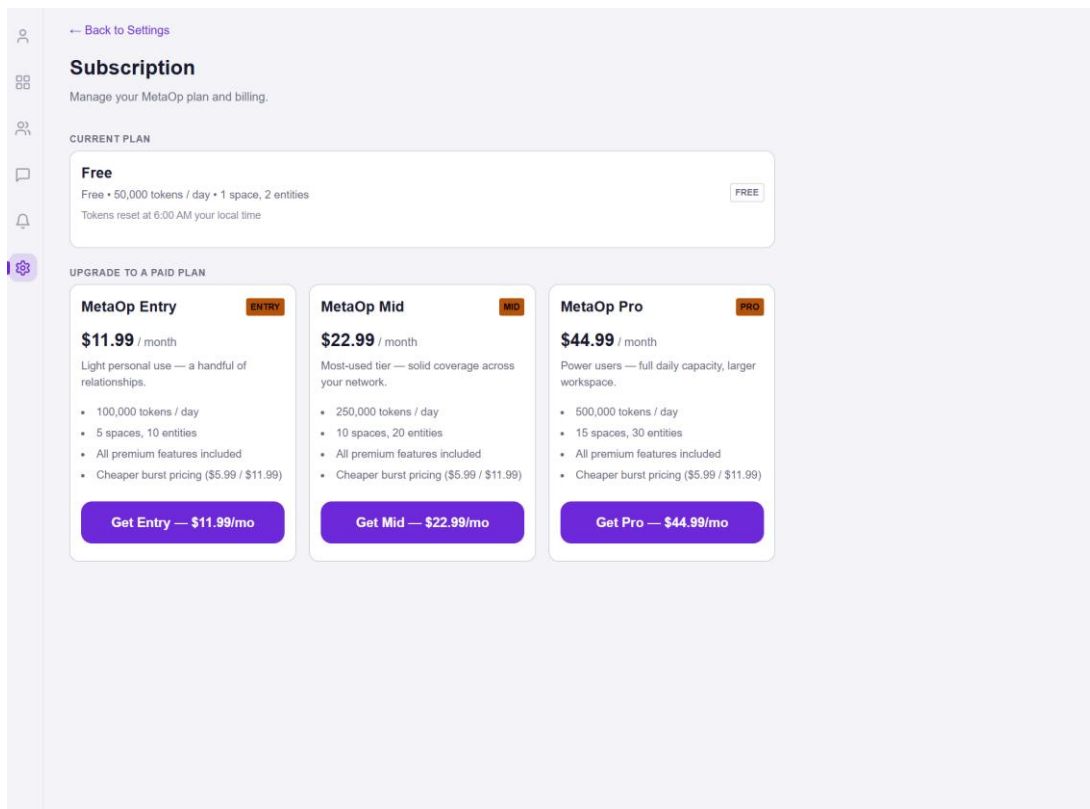


Figure AA, Depicts the subscription models clearly breaks down the tokens cost and subscription tiers where tokens are incentive but so are features and accessing to more features for heavier users.

Burst purchases - the escape hatch

When a user hits their daily cap and wants to continue, they can buy a Burst, a one-time token allotment that does not affect their subscription. Bursts are priced predictably and offered in two sizes (100K and 200K tokens), which expire 4 hours after purchase, shipped through the same store on web, iOS, and Android. The point is that the user is in control of their cost; the system is never surprising them with a bill, and the AI does not silently degrade when it runs low. Three

tiers, transparent caps, no hidden bloat. We charge for clarity, not for tokens. Burst tokens are available to both account types: Free and subscription.

Free accounts operate on a pay-as-you-go model for extra capacity. One-time token packs are available at \$11.99 for 100K tokens and \$29.99 for 200K tokens (these are one-time top-ups, not a subscription). Alternatively, the user can wait until 6 AM local time the next day, when the free daily allowance refreshes to 50K tokens at no cost to the user. We monetize impatience, not your data.

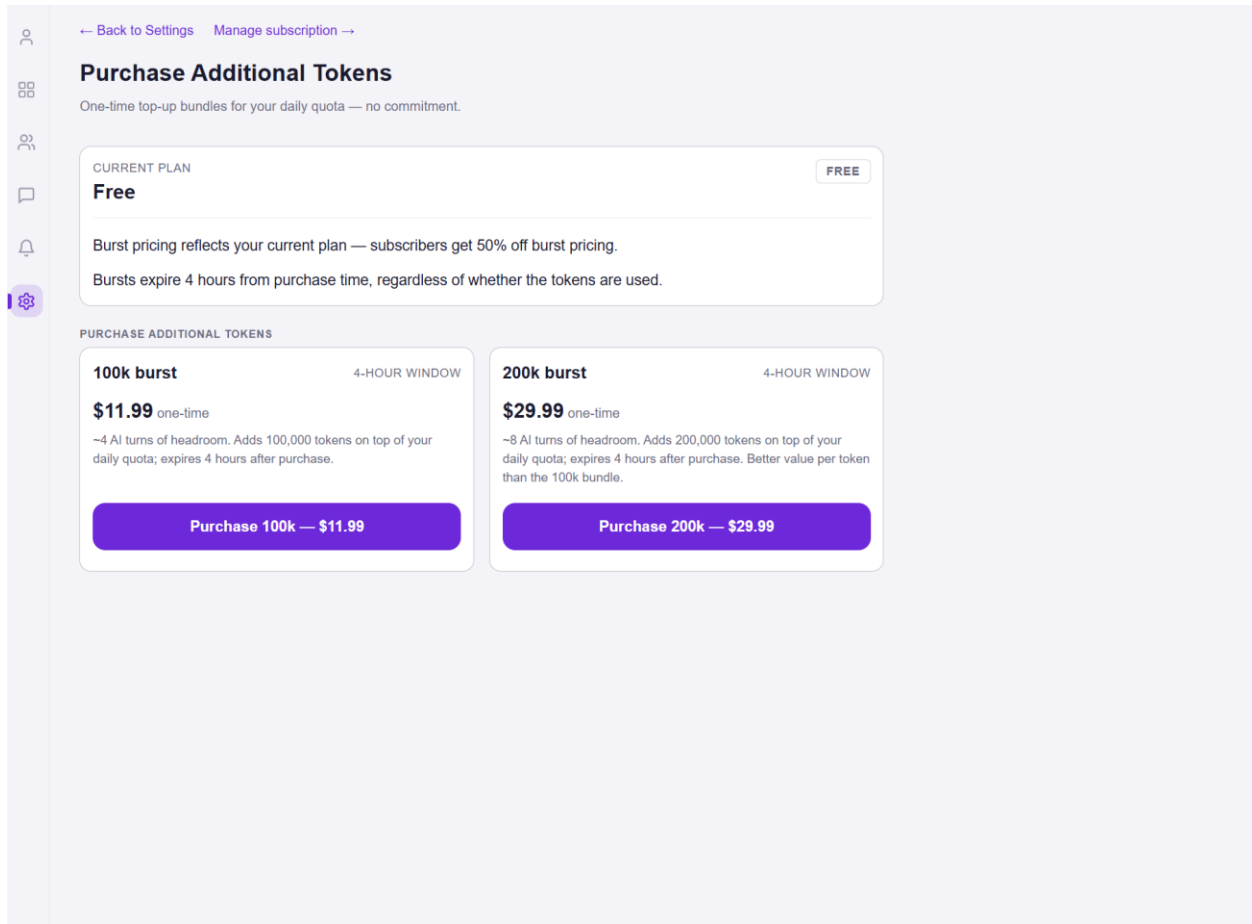


Figure BB, Additional tokens for free accounts pricing and limitations.

Subscribers who need additional tokens can purchase one-time burst packs, which expire 4 hours after purchase. As a subscriber benefit these are discounted 50%: \$5.99 for an additional 100K tokens and \$11.99 for an additional 200K tokens.

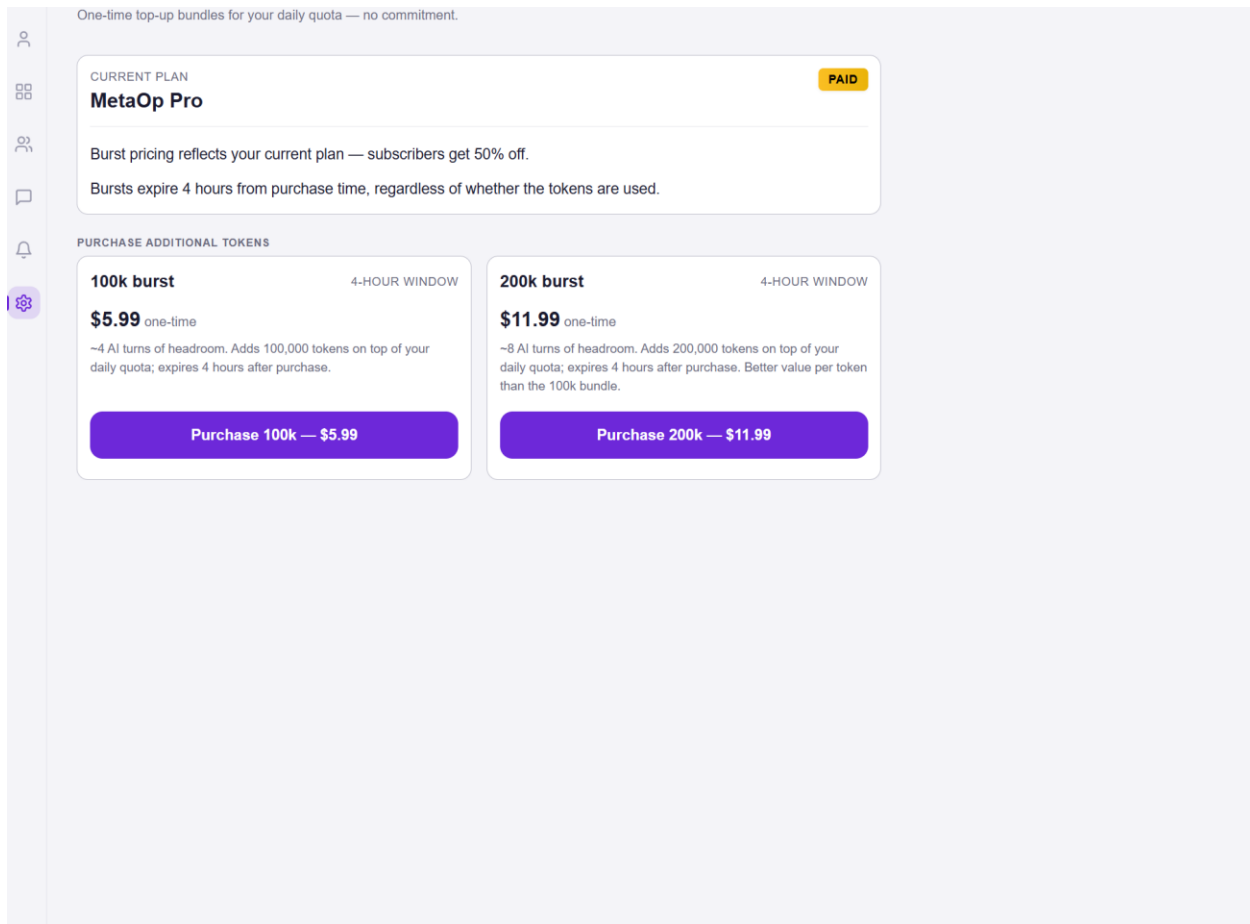


Figure CC, Additional tokens for subscription accounts pricing and limitations.

Why we price the way we do

Most AI is cheap or free because you're the product — your data monetized, your attention sold to advertisers. We don't do either. No data sales, no ads.

That means you pay for MetaOpAI directly. But because the architecture is built to be selective — extracting signal instead of replaying bloated context — your tokens go further. You're not paying to re-process the same conversation over and over; you're paying for precise, persistent memory. A premium from power users is designed to sustain free accounts for everyone else.

You pay a premium because, unlike most AI, we're not monetizing you.

Unit economics

We expect healthy gross margins because the substrate does most of the work and the LLM is doing rendering at the edge. A typical turn touches our substrate plenty but sends a small,

dense prompt to the model. Compared to a long-running wrapper conversation where the user is re-paying for the full transcript every turn, our per-turn LLM cost is lower by design. Where our costs are higher is in the infrastructure underneath, Postgres, Redis, async workers, observability, the trace store. Those are fixed costs that amortize across many users.

Path to platform economics

The consumer subscription is the wedge. The platform thesis is the longer arc. The substrate we built for personal cognition generalizes, a team's relational dynamics, a community's organizational drift, a family's intergenerational patterns. We are building the consumer product first because that is where the wedge is. The substrate stays open for the platform shape that emerges.

Defensibility and Moat

We think the moat is the substrate, not the model. The model is a rented commodity at the edge. The substrate is the durable asset, and the discipline behind it is what makes the substrate hard to copy.

The substrate compounds

Every cognition turn adds to the substrate. Every pattern that fires sharpens the calibration. Every trace that gets analyzed teaches us something. The substrate gets more valuable with time. A competitor that ships tomorrow would be at month zero of substrate maturity; we are at month many. That gap widens monotonically as long as the substrate keeps writing.

The discipline is harder to copy than the architecture

A competitor could read this paper and replicate the architecture diagram. What they could not replicate without years of iteration is the discipline encoded in the calibrations, the per-extractor confidence floors, the per-template density thresholds, the per-pattern decay coefficients, the per-scope retrieval shapes, the closed pattern vocabulary tuned against real narrations. Those values exist in our config files, not on the architecture slide.

Per-user substrate is the user's own moat

A user with twelve months of substrate would lose continuity if they switched products. That switching cost is real, and importantly, it belongs to the user. We make data export the first principle of the product so the user can take their continuity with them if they choose. But the continuity is built inside our substrate. Recreating it elsewhere requires re-narrating, re-extracting, re-pattern-detecting from scratch. Most users will not do that for the same reason most users do not move their email.

Distribution and the trust premium

We are competing for relational time, the moments when a person turns inward to think about their life. The product wins that competition on trust, on craft, and on the felt sense of being seen accurately. A wrapper that just bolted a memory feature onto a general assistant cannot match the trust premium of a product built for the relational domain from day one.

Traction and Validation

We are pre-launch. We do not have user counts to show, retention curves to plot, or revenue to disclose. What we have, instead, is technical validation that the substrate works in real production traffic, and architectural validation that the system can be evolved safely.

Technical validation

The substrate is shipped and running. The pattern engine has fired real patterns against real narrations. The cognition orchestrator runs the full intent-classification-and-routing path on every turn. The trace layer has caught and informed corrections to several non-trivial bugs (documented in the Trace Layer chapter). Migrations are clean. The QA harness validates the cognition pipeline end-to-end against four tier-appropriate test accounts.

Architectural validation

We have swapped the LLM provider mid-development without breaking downstream behavior, proof that the substrate is genuinely separable from the model. We have restructured the pattern engine's scoring functions and the orchestrator's policy resolution without rewriting the other layers, proof that the layer boundaries are real. We have introduced new pattern templates without touching the orchestrator, proof that the pattern engine is extensible by design.

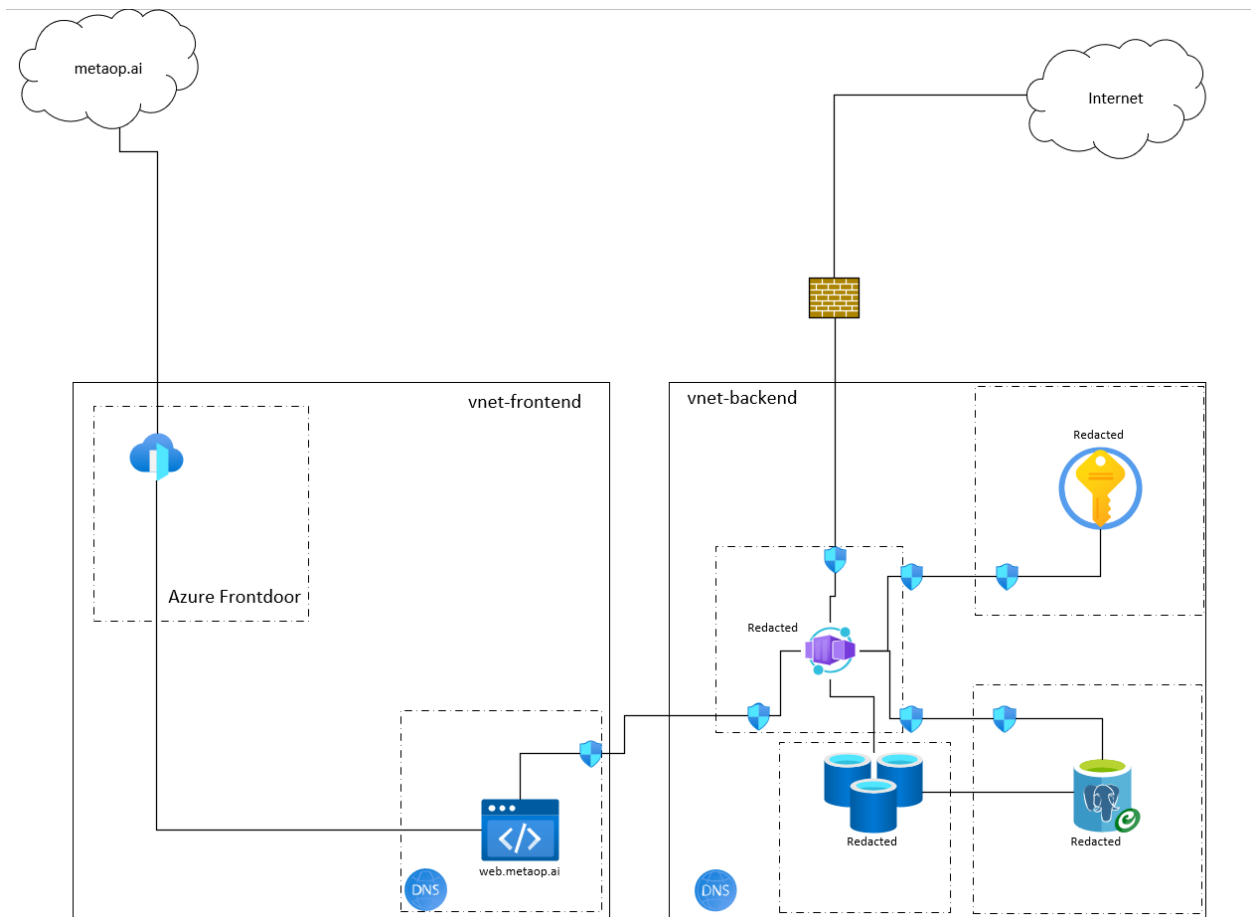
What we are validating next

Closed beta on production infrastructure. We are validating four things in parallel: latency under real concurrency, cost-per-turn under Azure OpenAI pricing with prompt caching enabled, pattern firing rate against real human narration, and qualitative felt-sense of continuity from invited beta users. The first three are instrumented through the trace layer. The fourth is interview-based and qualitative, and arguably the most important.

PART VIII

Platform

Deployed within Azure



[FIX the Figure DD caption] It currently reuses Figure CC's caption ("Additional tokens for subscription accounts pricing and limitations."), but Figure DD sits under "Deployed within Azure." Replace its caption with something like: "Figure DD, MetaOpAI's production deployment topology on Microsoft Azure." (adjust to the actual image)., Additional tokens for subscription accounts pricing and limitations.

Infrastructure and Security

MetaOpAI is built on a FastAPI backend using PostgreSQL for persistent storage and Redis for caching. The platform is designed to run in Azure Kubernetes Service (AKS) and uses Azure Managed Identity to securely access services such as Azure Key Vault, PostgreSQL, and Redis without storing credentials in code.

The architecture is segmented across multiple VNets with Network Security Groups (NSGs) enforcing micro-segmentation and controlling east-west (service-to-service) traffic, with north-south (ingress/egress) traffic governed at the perimeter. The frontend and backend are isolated into separate network segments, with VNet peering used for controlled communication between them.

The AKS environment is privately networked using private endpoints, meaning cluster management is abstracted through a secured local administration virtual machine that acts as the management plane access point for engineers and administrators.

The frontend application is built with React and Vite and is protected behind Azure Front Door with Web Application Firewall (WAF) enabled. Because much of the infrastructure relies on private endpoints, private DNS zones are required to properly resolve internal Azure services.

Outbound traffic initially routes through standard route tables and default routes without centralized firewall inspection. As the platform scales, firewall-based egress inspection and advanced routing controls can be introduced to provide additional security visibility and traffic governance.

Security

Data Security and Infrastructure Design

- **Encryption vs. Anonymization:** To ensure production-level utility, data is fully encrypted rather than anonymized. True anonymization often implies secondary data retention and commercialization; this architecture rejects that model in favor of strict, localized encryption at rest and in transit.
- **Identity and Access Management (IAM):** The platform completely offloads credential risk by enforcing exclusive federated authentication through trusted third-party providers (Microsoft, Google, and Apple).
- **Network Topology (Production):** Production environments utilize a secure, zero-trust routing model. Network Security Groups (NSGs) enforce strict north-south traffic segmentation, directing all backend routing through an enterprise firewall and a Web Application Firewall (WAF).
- **Architecture:** We employ private endpoints, so the majority of network traffic is routed over the Azure backbone and is not exposed externally. Administrative access is

governed by Azure RBAC and Privileged Identity Management (PIM), with device-bound, token-secured authentication to a bastion host.

- **Beta Environment Strategy:** To optimize infrastructure costs during the validation phase, the beta environment relies entirely on core architectural isolation and the WAF, deferring the full enterprise network suite until market viability justifies scaling.

Roadmap

Immediate (Q4 - 2026): Closed beta on Azure production infrastructure. Apple App Store Small Business Program enrollment. Stripe live mode cutover. App Store and Play Store submissions. Open beta. General availability.

Medium-term (2027): Platform expansion. Text to Voice, Team-scoped cognition (relational dynamics across an organization). Community-scoped cognition (group drift over time). Cross-platform integrations, calendar, messaging, email, that turn passive context into substrate.

Capital strategy

We are raising selectively. We are not a high-burn product, the substrate runs on commodity infrastructure, the LLM costs are bounded by the architecture's selectivity, and most of our compute is the cheap kind. We want capital partners who understand the cognition-substrate thesis and who can be patient through the substrate's compounding curve, which is slow at first and accelerates after the first six to twelve months of user data accumulate.

We are deliberately under-marketed today. The whitepaper is published; the website at <https://metaop.ai> is live; the contact at founder@metaop.ai is open. We prefer careful conversations to wide funnels. If this thesis resonates, the contact channel is real.

PART IX

Innovations

Innovations without constraints

What we Innovated

A wrapper is basically an AI chat app with a UI and prompts layered on top of a language model. It sends your conversation history to the model every turn and treats the conversation itself as memory.

MetaOpAI tries to build something deeper: a system with its own structured memory and cognition layer underneath the model.

Here's the simplified breakdown:

- **Structured memory (KRL):**
Instead of one giant chat history, the system stores information in organized categories:
 - you,
 - another person,
 - an environment,
 - or a relationship.Each record is searchable, traceable, and confidence-scored.
- **Extraction instead of summaries:**
Rather than compressing chats into vague summaries, the system extracts specific signals, events, and patterns with evidence attached.
- **Noise filtering:**
Multiple filters remove weak, irrelevant, or outdated information so only the most useful context reaches the AI.
- **Pattern recognition:**
The system looks for recurring behaviors and dynamics across time, not just what happened in the latest message.
- **Time weighting:**
Recent events matter more than old ones, but important older events are still retained instead of forgotten.
- **Different cognition modes:**
The system can think:
 - about you personally,
 - about a relationship,
 - or about a larger environment like work or family.Each uses a different reasoning style.
- **Configurable behavior:**
Retrieval depth, tone, and consistency can be adjusted without rewriting the system.

- **Traceability:**
Every conclusion can be traced back to where it came from and which extractor created it.
- **Async cognition updates:**
The memory system can update in the background while keeping the experience fast.
- **Cold-start scaffolding:**
Temporary AI-generated context can help early conversations, but it is blocked from becoming permanent memory unless supported by real user narration.
- **Temporal continuity:**
The system tries to understand story arcs over time:
 - escalation,
 - withdrawal,
 - rebuilding trust,
 - recurring cycles,
rather than isolated moments.
- **Model independence:**
The LLM itself is replaceable.
The real value lives in the continuity substrate, not the model provider.

The core idea is simple:

A wrapper replays conversations.

MetaOpAI tries to build a structured cognition system that compounds understanding over time.

PART X

Lessons Learned

What we've learned and Challenges

What we've learned

Lessons That Shaped MetaOpAI

The architecture did not emerge fully formed.

Most of the important design decisions came from observing failures in real runtime conditions and correcting them.

1. Selective Intelligence

The system became better when we stopped trying to make it do more and started teaching it when to do less.

Early versions ran maximum cognition on every message:

- deep pattern analysis,
- citations,
- retrieval,
- continuity injection.

The result felt unnatural.

A simple message like:

“Feeling alright today”

would receive an overly analytical response.

So, we introduced **Selective Intelligence**:
the orchestrator now decides whether a message deserves:

- lightweight conversation,
- or deep analytical cognition.

Most turns intentionally stay shallow.

Lesson

Restraint is part of what makes cognition systems feel human.

2. Epistemic Integrity

We learned the system must be disciplined about what it actually knows.

Several problems forced architectural corrections:

Ontology Drift

The model sometimes invented pattern categories that did not exist.

We fixed this with ontology guards so the system can only surface patterns that truly exist in the substrate.

Recursive Contamination

If a memory system ingests its own AI-generated interpretations as truth, it slowly drifts into hallucination loops.

We solved this by separating:

- synthetic cognition,
- from durable cognition.

AI-generated scaffolding can never become permanent evidence.

Cold-Start Emptiness

New users had no continuity, so the system initially felt hollow.

We introduced a temporary low-authority bootstrap layer (“hourglass cognition”) that fades as real continuity accumulates.

Lesson

A cognition system must be more careful about what it refuses to claim than what it asserts.

3. Extraction vs Summarization

Most AI systems repeatedly summarize growing conversations.

That compression is lossy:

- nuance disappears,
- timing disappears,
- confidence disappears,
- relational context disappears.

MetaOpAI treats this as a systems problem.

Instead of carrying transcripts forever, the system extracts:

- signals,
- events,
- patterns,
- context,
- provenance,
- and confidence.

Lesson

Summarization compresses and forgets.

Extraction structures and preserves.

4. Preventing Feedback Loops

A persistent AI system can slowly start believing its own interpretations more than the user's real narrative.

Over time:

- summaries become context,
- context becomes future interpretation,
- and abstraction replaces reality.

We treated this as one of the most dangerous risks in long-term memory systems.

So, we enforced strict separation between:

- the read path,
- and the write path.

Only user narration becomes durable cognition.

Lesson

A long-term cognition system must never trust its own summaries more than the user's lived experience.

5. Observability and Debugging

Most major corrections came from cognition traces, not intuition.

Examples:

- stream hangs caused by enrichment bottlenecks,
- users repeating questions before cognition settled,
- async extractors causing stale state.

This led to:

- structured tracing,
- consistency modes,
- settle markers,
- and silent UI refreshes.

Lesson

If a cognition system cannot explain what it did, debugging eventually collapses into interrogating the model itself.

6. Layered Architecture

The earliest runtime was heavily coupled and difficult to evolve.

Separating the system into layers changed that:

- journal,
- substrate,
- pattern engine,
- orchestrator,
- composer.

This made it possible to:

- swap providers,
- evolve retrieval,
- change scoring,
- and add cognition features without destabilizing the whole system.

Lesson

Layering makes cognition systems evolvable.

7. Reality Corrected the Architecture

Many important invariants were not part of the original design:

- ontology guards,
- provenance enforcement,
- trace layers,
- temporal continuity,
- cognition consistency,
- synthetic cognition separation.

They emerged because runtime behavior exposed weaknesses.

Final Lesson

The architecture matured by repeatedly observing failures, tracing their causes, and allowing reality to correct the system.

BACK MATTER

Honest Limits, Glossary, Contact

What we do not know yet, the terms we used along the way, how to reach us.

What We Do Not Know Yet

The honest list. Open questions, unproven hypotheses, things that might be wrong.

We do not yet know how the cognition layer behaves at scale across thousands of users with months of accumulated history. We have run the substrate against synthetic and small-scale real traffic, and the architectural assumptions have held. We expect at scale to discover failure modes we have not seen. Some will be performance-shaped. Some will be epistemic, patterns that fire too aggressively, or not aggressively enough, on cohorts we have not yet seen.

We do not know whether the three-scope model (USER / ENTITY / SPACE) is the final shape or a way station. Some of the most interesting cognition we have seen lives in the intersection between scopes, and we may discover that a fourth or fifth scope is warranted. We have left room for it.

We do not know whether the right business model is subscription, usage-metered, freemium-with-burst, or something else. We have shipped a three-tier subscription with usage caps and burst purchases for those who want more. We will let real traffic teach us what fits.

We do not know how the broader AI industry will price long-term memory as it commoditizes. Some of what we built may end up being a default feature of the major model providers in eighteen months. We think the substrate-level discipline, the filters, the scopes, the trace layer, will still matter even if the underlying memory primitive becomes a commodity. We could be wrong about that.

We do not know how much the pattern engine's pattern templates will need to evolve as we see more domains. The current twenty-six pattern templates (distinct from the twenty substrate cells of the KRL matrix) were tuned against relationships and personal growth. Other domains, work performance, health, learning, will likely need their own template families. We are building the bench, not just the templates.

We are not certain. We are committed. The certainty grows as the traces grow.

Glossary

Cognition substrate

The persistent, structured layer between the LLM and the user experience. Stores signals, events, meta-context, context, and patterns across four subject scopes.

KRL (Knowledge Representation Layer)

The cognition substrate's matrix of four subject scopes (USER, ENTITY, SPACE, RELATIONSHIP_PAIR) by five-layer kinds (Signal, Event, Meta-Context, Context, Pattern). Twenty addressable cells per user.

Signal

A per-utterance behavioral marker extracted from a single narration. Example: withdrawal, kept_commitment, repair_attempt_made.

Event

A discrete moment narrated by the user, optionally with a temporal anchor. Example: argument with brother on May 17.

Meta-Context

A drift-prone interpretive read, the user's own framing of what is happening, separable from the underlying factual context. Example: "I think the relationship is getting better."

Context

A stable factual fact about the user or an entity. Example: user works in healthcare; user's sister moved abroad.

Pattern

A typed artifact that fires when a configurable signal aggregate crosses a threshold across signals that span days or weeks. Carries confidence, evidence, awareness tier, action tier, and lineage.

Cognition orchestrator

The runtime traffic controller. For each turn, classifies intent, resolves scope, fetches continuity, and assembles the directive block the composer will use.

Cognition policy

A small JSON contract that tunes orchestrator behavior per scope (retrieval_depth, citation_mode, consistency_mode, etc.). Resolves through a four-tier hierarchy.

Composer

The layer that constructs the LLM prompt with the right context, citations, and safety overlays. Returns the user-visible response.

Provenance

The JSONB column on every extracted record that names the extractor, version, source message, and timestamp. A Phase L-1 invariant.

Trace

A structured record of one cognition turn, what the orchestrator decided, what was fetched and why, what patterns fired, what the composer assembled, what the LLM said. The audit trail.

Synthetic cognition

A session-scoped JSONB blob that holds low-authority bootstrap material during cold-start. Tapers as durable cognition accumulates. The hourglass.

Arc detector

A pattern detector that fires when a recognizable temporal sequence of signals occurs, RepairThenRelapse, TrustRebuilding, EscalationArc, WithdrawalSpiral.

Selective intelligence

The discipline that the system does less, more accurately, more often. Operationalized as the three-level cognition hierarchy.

Be Understood

The motto. Also, the product strategy. The deepest human need an AI can serve.

Contact and Closing

If this paper resonated with you, whether you are an engineer, an investor, a potential collaborator, or someone who has felt the cognition gap in their own AI tools and wants to talk about what we are building, we would rather have a careful conversation than a wide spray of pitches.

Website: <https://metaop.ai>

Founder: founder@metaop.ai

Founder: Tony S.

We did not have all the answers when we started. We do not have them now. What we have is a substrate, a discipline, and a thesis we have been willing to be corrected by. If that is your kind of conversation, we would like to have it.

*MetaOpAI Whitepaper, Unified
May 2026*