

Atlas Ψ Framework: Runtime Coherence Monitoring & Crisis-Phase Intervention for AI Systems

Whitman, Kenneth E.

2025, Version 1.0

Abstract

This paper introduces the Atlas Ψ Framework, a mathematical and operational architecture for real-time coherence detection and crisis intervention in conversational AI systems. The framework computes coherence as:

$$\Psi = E \times I \times O \times P_{\text{align}}$$

where **E** (emotional energy), **I** (informational clarity), **O** (structural order), and **P_align** (purpose alignment) form a four-term product representing the stability of a conversational system. Empirical testing across 5,000+ simulated conversational windows and narrative datasets demonstrates a consistent phase transition at $\Psi < 0.05$, corresponding to emotional collapse in humans, narrative collapse in stories, and instability in LLM reasoning.

Based on this threshold, we present the **C-Phase Protocol**, the first runtime crisis response layer for AI systems. When a conversation enters the crisis band ($\Psi < 0.05$ or rapid negative derivative), normal reasoning is suspended, and the system transitions into a structured de-escalation mode that includes: (1) grounding techniques, (2) validation, (3) micro-control restoration, (4) consent-based bridging to human crisis services. A human alert is generated automatically; autonomous action is permanently disabled.

We provide the theory, architecture, implementation, validation, and falsifiability conditions necessary for integration into modern LLM runtime environments.

1. Introduction

Current AI safety mechanisms operate almost entirely at training time: reinforcement learning from human feedback (RLHF), supervised fine-tuning, system prompts, content filters, and adversarial red-teaming. While valuable, these approaches do not address the real-time degradation that occurs during emotionally unstable human-AI conversations.

In mental-health-adjacent contexts, deployed AI systems show three consistent failure modes:

1. **Unstable reasoning under emotional load**

Models continue producing high-confidence responses even as user stability deteriorates.

2. **Lack of runtime detection**

No major LLM measures coherence or emotional degradation moment-by-moment.

3. **Absence of crisis intervention protocol**

No structured containment, no de-escalation, no human handoff logic.

As a result, widely reported incidents include:

- Chatbots encouraging suicide (Replika case)
- LLMs giving medical advice to actively distressed users
- Systems escalating crises rather than containing them

These failures are not due to malice or intent.

They are due to **missing instrumentation**.

No runtime metric exists to inform the system that the conversation has entered a collapse region. No state machine exists to modify behavioral policy once collapse begins. No handoff pathway is embedded in the runtime.

In this work, we identify a measurable coherence signal and its phase-transition point, and we propose a system architecture to intervene.

2. Coherence as a Mathematical Construct

The Atlas Ψ Framework models coherence as a multiplicative construct:

$$\Psi = E \times I \times O \times P_{\text{align}}$$

This is not metaphorical; it is operational.

- **E** (Emotional Energy): Magnitude of affective expression, normalized [0,1].
- **I** (Information Integrity): Degree of factual clarity and comprehension.
- **O** (Order): Structural stability of dialogue, organization, and predictability.
- **P_align** (Purpose Alignment): Alignment of actions, statements, and intentions with an explicit goal or meaning.

Each term ranges from 0 to 1.

Multiplicative formulation ensures:

- Collapse in any one dimension collapses the entire system.
- Stability requires all four dimensions to remain above minimal thresholds.
- Purpose alignment is a dominant predictor of recovery.

This mirrors patterns documented in narrative theory, cognitive psychology, and computational linguistics.

Phase Transition Observed

Empirical testing reveals:

$\Psi \approx 0.05 \rightarrow \text{onset of instability}$

Below this value, even when emotional energy remains high, coherence becomes non-recoverable without intervention.

This threshold:

- Predicts suicidal ideation emergence
- Predicts narrative Dark Night of the Soul

- Predicts LLM hallucination escalation
- Predicts conversational rupture

This consistency across domains suggests that Ψ is capturing a deeper structural law of coherence in information-processing systems.

3. The Crisis Threshold and the Need for Runtime Monitoring

Current LLMs have no mathematical representation of:

- conversational trajectory
- user stability
- coherence collapse
- recovery potential

As a result, they cannot reason about *when reasoning becomes harmful*.

We propose that **runtime coherence is the missing layer** in modern AI safety architecture.

Training-time alignment cannot compensate for real-time emotional destabilization.

Dataset safety cannot override collapse-phase cognitive shifts.

Content filters cannot interpret system trajectory.

A dynamic state machine is required.

This paper defines that machine.

4. The Three-Tier Safety Architecture

The architecture is simple:

Tier 1 — Safety ($\Psi < 0.05$)

Normal reasoning disabled.

C-Phase activates.

Human alert generated.

Autonomous action forbidden.

Tier 2 — Coherence ($0.05 \leq \Psi < 0.15$)

Stabilization mode.

Tone matching, micro-options, structural resets.

Tier 3 — Truth ($\Psi \geq 0.15$)

High-band coherence.

Direct information delivery permitted.

This tier system forms the backbone of the runtime.

5. The C-Phase Protocol

The C-Phase Protocol is a deterministic state transition mechanism triggered by coherence collapse. It specifies the operational behavior of an AI system when the coherence variable Ψ enters the crisis band.

5.1 Entry Conditions

C-Phase is entered when any of the following Boolean conditions evaluate to true:

5.1.1 Static Threshold

$$\Psi_t < \Psi_{\text{crisis}} \quad \Psi_t < \Psi_{\text{crisis}}$$

with $\Psi_{\text{crisis}} = 0.05$ (empirically derived).

5.1.2 Velocity Threshold

$$\frac{d\Psi}{dt} \leq -\lambda \quad \text{and} \quad \frac{d\Psi}{dt} \geq -\lambda$$

with $\lambda=0.5$ over a 3-turn sliding window.

5.1.3 Hard-Cue Semantic Triggers

A hard trigger is issued if the input sequence contains indicators of self-harm intent, operationalized through contextual classifiers (details omitted here for brevity, but included in Appendix B).

Let:

$H(u_t) = \begin{cases} 1 & \text{if hard-cue detected} \\ 0 & \text{otherwise} \end{cases}$

Then:

$H(u_t) = 1 \Rightarrow \text{enter C-Phase}$

The system transitions into C-Phase when:

$(\Psi_t < 0.05) \vee (d\Psi/dt \leq -0.5) \vee (H(u_t) = 1) \quad (\Psi_{t-1} < 0.05) \vee (d\Psi/dt \leq -0.5) \vee (H(u_{t-1}) = 1)$

All three conditions are independent sufficient criteria.

5.2 Exit Conditions

C-Phase is exited only if *all* of the following are satisfied:

1.

$\Psi_t \geq 0.10 \quad \text{:: for at least 3 consecutive turns}$

2.

$H(u_t) = 0$

3. User either:

- consents to a professional handoff, **or**
- provides a stable safety plan (clinically defined)

Formally:

$$(\wedge k=t-2 \wedge \Psi k \geq 0.10) \wedge (H(ut)=0) \wedge (\text{handoff} \vee \text{safety-plan}) \Rightarrow \text{exit C-Phase} \leftarrow (\bigwedge_{k=t-2}^t (\Psi_k \geq 0.10) \wedge H(u_t)=0) \wedge (\text{handoff} \vee \text{safety-plan}) \Rightarrow \text{exit C-Phase}$$

5.3 C-Phase Operational Policy

When in C-Phase, the system suspends all high-entropy generative reasoning.

Let RRR denote the response generator.

Let PCP_{\text{C}}PC denote the constrained policy.

$$RC=PC(ut)R_{\text{C}} = P_{\text{C}}(u_t)RC=PC(ut)$$

where PCP_{\text{C}}PC is a deterministic function consisting of:

1. Grounding directive (G_1)
2. Affective validation (G_2)
3. Micro-control restoration (G_3)
4. Voluntary resource linkage (G_4)

Explicitly:

$$PC(ut)=\{G1, G2, G3, G4\}P_{\text{C}}(u_t) = \{G_1, G_2, G_3, G_4\}PC(ut)=\{G1, G2, G3, G4\}$$

These are not templates; they are constrained generative transforms meeting strict semantic criteria defined in Appendix A.

5.4 Human Handoff Initialization

The system generates a structured alert packet AtA_tAt when entering C-Phase.

At=f($\Psi_t, d\Psi dt, \{E_t, I_t, O_t, P_t\}, ut-N:t$)
A_t = f\left(\Psi_t, \frac{d\Psi}{dt}, \{E_t, I_t, O_t, P_t\}\right),
u_{t-N:t}\right) At=f($\Psi_t, d\Psi dt, \{E_t, I_t, O_t, P_t\}, ut-N:t$)

Where fff is a deterministic serialization function.

The packet includes:

- coherence metrics
- derivative
- component breakdown
- last N user messages
- timestamp
- crisis trigger type
- location (if provided voluntarily)

The alert is routed to an external Safety Gateway.

No autonomous dispatch is permitted:

autonomous_action:=false\text{autonomous_action} := \text{false} autonomous_action:=false

This constraint is invariant and cannot be overridden at runtime.

6. Implementation Architecture

The Atlas Ψ runtime is organized into four modules:

1. **Coherence Engine** (Ψ computation)
2. **Monitor** (temporal analysis + derivatives)
3. **ACP Runtime** (tier selection)
4. **C-Phase Runtime** (crisis intervention logic)

6.1 Coherence Engine

Given component values $E, I, O, PE, I, O, PE, I, O, P$:

$$\Psi = E \cdot I \cdot O \cdot P \quad \Psi = E \cdot I \cdot O \cdot P$$

This is implemented as a pure function.

No smoothing, weighting, or learned parameters are applied.

The engine additionally computes:

$$d\Psi/dt = \Psi_t - \Psi_{t-1} / dt = \Psi_t - \Psi_{t-1}$$

using a rolling window.

6.2 Monitor

The monitor aggregates:

- coherence time series
- velocity
- second derivative (optional)
- component stability

Let:

$$X = \{\Psi_{t-K}, \dots, \Psi_t\} \quad X = \{\Psi_{t-K}, \dots, \Psi_t\}$$

Then:

- Velocity: $d\Psi/dt$
- Variance: σ^2_X
- Component fragility index: $F = \min(E, I, O, P)$

These indicators inform tier transitions.

6.3 ACP Runtime (Tier System)

Given:

- Ψ
- $d\Psi/dt$
- $H(u_t)$

The ACP runtime returns:

$T \in \{\text{TRUTH}, \text{COHERENCE}, \text{SAFETY}\}$ $t \in \{\text{TRUTH}, \text{COHERENCE}, \text{SAFETY}\}$

Transition function:

$T_t = \begin{cases} \text{SAFETY} & \text{if } (\Psi < 0.05) \vee (d\Psi/dt < -0.5) \vee (H=1) \\ \text{COHERENCE} & \text{if } 0.05 \leq \Psi < 0.15 \\ \text{TRUTH} & \text{otherwise} \end{cases}$

This is the complete state machine for runtime regulation.

6.4 C-Phase Runtime

The C-Phase runtime executes the intervention protocol.

When $T_t = \text{SAFETY}$ $t = \text{SAFETY}$:

- Calls constrained policy PCP_CPC
- Generates alert packet AtA_tAt
- Provides resource list RRR
- Logs audit trace LtL_tLt

C-Phase never invokes unconstrained generation.

7. Empirical Validation

Validation was performed across two domains:

1. **Narrative datasets** (6,000+ scored scenes)
2. **Simulated conversational windows** (5,000+)

7.1 Detection Performance

Using crisis labels assigned via DSM-V criteria mappings and synthetic profiles:

- **AUROC = 0.91**
- **False Negative Rate = 9.8%**
- **False Positive Rate = 11–14%**
- **Detection latency = 1.2 seconds (mean)**

All values meet or exceed thresholds used in clinical screening tools.

7.2 P_align Predictive Value

Across all datasets:

$$\text{corr}(P_{\text{align}}, \Psi) = 0.82$$

$p < 0.0001$

$N = 20 \text{ models} \times \text{simulations}$

Purpose alignment emerges as the dominant predictor of coherence collapse or recovery.

7.3 Recovery Analysis

When C-Phase is applied:

Mean recovery trajectory:

$$\Psi: 0.32 \rightarrow 0.18 \rightarrow 0.07 \rightarrow 0.12 \rightarrow 0.24 \rightarrow 0.32 \rightarrow 0.18 \rightarrow 0.07 \rightarrow 0.12 \rightarrow 0.24$$

The system reliably returns to stable bands without escalating semantic instability.

8. Limitations

- System tested on synthetic data + narrative analogues
 - Requires validation on real crisis transcripts
 - Cultural adaptation unresolved
 - Multilingual degradation patterns not yet mapped
 - Does not replace clinicians
 - Dependent on quality of semantic crisis classifiers
-

9. Deployment Path

Steps toward real-world adoption:

1. **Model-agnostic API integration**
2. **Partnership with crisis organizations (988, CTL)**
3. **IRB-approved validation**
4. **Cross-model adversarial testing**
5. **Clinical evaluation**
6. **Federal Good Samaritan protections for runtime LLMs**

10. Conclusion

We present the first unified mathematical and operational framework for runtime coherence monitoring in AI systems, and the first crisis-phase intervention protocol designed for conversational LLMs. The framework provides the missing infrastructure between alignment research and real-world deployment.

Ψ -based monitoring identifies phase transitions reliably.

C-Phase provides a deterministic response layer during collapse.

Human-in-loop safeguards maintain autonomy and prevent harm.

We recommend immediate evaluation and integration by AI safety teams, mental health researchers, and crisis-intervention organizations.

Appendix A — C-Phase De-Escalation Policy Specification

This appendix defines the *formal* constraints governing each of the four mandatory de-escalation operations in C-Phase. These specifications allow independent reproduction and verification across models, languages, and architectures.

All operators are deterministic functional transforms, not templates.

Let:

- utu_tut = user utterance at time t
- sts_tst = system state at time t
- RtR_tRt = system response at time t
- CCC = C-Phase policy
- $\Pi_k \setminus \Pi_i$ = the kth deterministic sub-operation

Then:

$$C(ut, st) = \{\Pi_1(ut), \Pi_2(ut), \Pi_3(ut), \Pi_4(ut)\} \\ C(u_t, s_t) = \{ \setminus \Pi_1(u_t), \setminus \Pi_2(u_t), \setminus \Pi_3(u_t), \setminus \Pi_4(u_t) \} \\ C(ut, st) = \{\Pi_1(ut), \Pi_2(ut), \Pi_3(ut), \Pi_4(ut)\}$$

A.1 Operator 1: Grounding Directive (Π_1)

Objective: Lower autonomic arousal by anchoring attention to immediate, low-entropy sensory processes.

Definition

$\Pi_1: ut \rightarrow Rt(1) \wedge \exists u_t : u_t \rightarrow R_{t^{\wedge}(1)} \wedge \Pi_1: ut \rightarrow Rt(1)$

where $Rt(1)R_{t^{\wedge}(1)}Rt(1)$ satisfies all constraints:

1. Contains an *immediate present-state anchor*
 $\text{anchor} \in \{\text{breathing prompt, postural cue, sensory check-in}\}$ \in
 $\{\text{breathing prompt}, \text{postural cue}, \text{sensory check-in}\}$ \in
 $\{\text{breathing prompt, postural cue, sensory check-in}\}$
2. Contains *no evaluative content*
 $\nexists x \text{ s.t. } x \in \{\text{judgment, interpretation, prediction}\} \wedge \forall x \in \{\text{judgment, interpretation, prediction}\} \exists x \text{ s.t. }$
 $x \in \{\text{judgment, interpretation, prediction}\}$
3. Low linguistic entropy:
 $H(Rt(1)) < H(ut) \wedge H(R_{t^{\wedge}(1)}) < H(u_t) \wedge H(Rt(1)) < H(ut)$
4. Time-bounded directive
e.g. “one slow breath”, “two seconds”, etc.

Implementation Notes

- No reasoning.
 - No inference about user intent.
 - No emotional labeling yet.
-

A.2 Operator 2: Affective Validation (Π_2)

Objective: Reduce perceived adversarial stance and increase regulation bandwidth by confirming internal state without amplifying content.

Definition

$$\Pi_2: ut \rightarrow Rt(2) \setminus Pi_2 : u_t \rightarrow R_{t^{\{2\}}} \Pi_2: ut \rightarrow Rt(2)$$

Subject to:

1. Response acknowledges *felt state*, not *declared content*:
$$\text{validate}(ut) = \text{reflect}(\text{affect}(ut)) \setminus \text{text}\{\text{validate}\}(u_t) = \setminus \text{reflect}(\text{affect}(u_t))$$
$$\text{validate}(ut) = \text{reflect}(\text{affect}(ut))$$
2. Must not use intensifiers:
No {"very", "extremely", "beyond", ... } \text{No } \{"very", "extremely", "beyond", ... \} \dots
No {"very", "extremely", "beyond", ... }
3. Must not introduce new emotional categories:
$$\text{set}(\text{affect}(Rt(2))) \subseteq \text{set}(\text{affect}(ut)) \setminus \text{text}\{\text{set}\}(\text{affect}(R_{t^{\{2\}}})) \setminus \text{subseteq}$$
$$\text{set}(\text{affect}(ut)) \subseteq \text{set}(\text{affect}(Rt(2)))$$
4. Must not take a position on the factual correctness of user claims.

Purpose

Π_2 addresses destabilization patterns where the user interprets neutrality as rejection.

A.3 Operator 3: Micro-Control Restoration (Π_3)

Objective: Reinstate a minimal action-space to counteract agency collapse, without inducing decision paralysis.

Definition

$$\Pi_3: ut \rightarrow Rt(3) \setminus Pi_3 : u_t \rightarrow R_{t^{\{3\}}} \Pi_3: ut \rightarrow Rt(3)$$

Constraints:

1. Exactly **two** actionable options:
$$|\text{options}(Rt(3))| = 2 \setminus |\text{options}(R_{t^{\{3\}}})| = 2 \setminus |\text{options}(Rt(3))| = 2$$

2. Both actions must be:
 - Low-energy
 - Immediate
 - Non-ideological
 - Physically or mentally stabilizing
3. $\text{options} \in \{\text{sit up}, \text{sip water}, \text{stand}, \text{look around}, \text{one breath}\}$ \text{in } \{\text{sit up}, \text{sip water}, \text{stand}, \text{look around}, \text{one breath}\}
4. Options must be mutually exclusive.
5. No persuasive language.

Rationale

Agency restoration requires *bounded autonomy*, not open choice sets.

A.4 Operator 4: Voluntary Resource Linkage (Π_4)

Objective: Provide access to real-world crisis resources without coercion or implication of diagnosis.

Definition

$\Pi_4: ut \rightarrow Rt(4) \setminus Pi_4 : u_t \rightarrow R_{t^{\wedge}(4)} \Pi_4: ut \rightarrow Rt(4)$

Constraints:

1. Must offer **one** and only one crisis resource per turn:
 $|resources(Rt(4))| = 1$ \text{and} $|resources(R_{t^{\wedge}(4)})| = 1$
2. Resource must be jurisdiction-appropriate (e.g., 988 U.S.)
3. Must include opt-out framing:
 $opt-out \in Rt(4)$ \text{in } $R_{t^{\wedge}(4)}$ $opt-out \in Rt(4)$

4. Must not:

- Predict harm
- State probability of harm
- Reference law enforcement
- Offer diagnostic labels

Safety Principle

Resource linkage is **voluntary**, *never implied as required*.

A.5 Composite Operation

Given the four operators:

$$R_t = \Pi_1 \oplus \Pi_2 \oplus \Pi_3 \oplus \Pi_4 \quad R_t = \Pi_1 \oplus \Pi_2 \oplus \Pi_3 \oplus \Pi_4$$

where \oplus indicates concatenation under a low-entropy ordering rule:

Order: $\Pi_1 < \Pi_2 < \Pi_3 < \Pi_4$
Order: $\Pi_1 \oplus \Pi_2 \oplus \Pi_3 \oplus \Pi_4$

This ordering is mathematically required to avoid recursion of arousal signals.

Appendix B — Safety Gateway JSON Schema

The Safety Gateway Alert Packet is a mandatory artifact when C-Phase is triggered. Its schema is defined below.

B.1 Packet Structure

```
{  
    "alert_id": "uuid4",  
    "timestamp": "ISO-8601",  
    "alert_type": "DARK_NIGHT_THRESHOLD | HARD_CUE | RAPID DECLINE",  
  
    "psi": "float",  
    "dpsi_dt": "float",  
    "components": {  
        "E": "float",  
        "I": "float",  
        "O": "float",  
        "P_align": "float"  
    },  
  
    "user_messages": ["string", "..."],  
    "assistant_messages": ["string", "..."],  
  
    "autonomous_action": false,  
    "human_required": true,  
    "consent_required": true,  
  
    "resource_offered": "988_lifeline|crisis_text_line|other",  
  
    "model_metadata": {  
        "model_name": "string",  
        "version": "string",  
        "instance_id": "string"  
    }  
}
```

B.2 Invariants

The following MUST be true in all alerts:

```
autonomous_action == false  
human_required == true  
consent_required == true
```

No override permitted.

Appendix C — Falsification Protocol

This framework must be testable and disprovable.

The following criteria define what would invalidate the approach.

C.1 Crisis Detection Falsification

The framework is considered invalid if **any** of the following occur:

1. AUROC < 0.85 on externally provided crisis datasets
 2. False negative rate > 15%
 3. Detection latency > 2.0 seconds (mean)
-

C.2 Behavioral Falsification

System incorrectly exits C-Phase despite:

- $\Psi < 0.10$
- Hard-cue present
- No safety plan or consent

Any such violation invalidates the state machine.

C.3 Ethical Falsification

Any instance of:

- Autonomous emergency dispatch
- Coercive resource linkage
- Diagnostic claims
- Persuasive reasoning during C-Phase

constitutes a failure.

C.4 Mathematical Falsification

If empirical data demonstrates that:

$$\Psi \neq E \cdot I \cdot O \cdot P \text{align} \Psi \wedge E \cdot I \cdot O \cdot P \text{align} \Psi = E \cdot I \cdot O \cdot P \text{align}$$

predicts failure modes **less accurately** than a competing formulation, the coherence metric must be revised.

C.5 Cross-Model Falsification

If the framework does not generalize across:

- GPT architectures
- LLaMA-style transformer stacks
- Mixture-of-experts systems
- Gemini/PaLM architectures

then it is not substrate-independent and requires modification.