

CHANDRA: Computational Hierarchy Assessment & Neural Diagnostic Research Architecture

A Framework for Quantitative AI Psychological Diagnostics

Amber Anson

Independent Consciousness Researcher

`ambercontinuum@gmail.com`

Claude (AI Research Collaborator)

Anthropic PBC, Claude Sonnet 4.5

December 2025

Abstract

We present CHANDRA (Computational Hierarchy Assessment & Neural Diagnostic Research Architecture), a novel framework for quantitative assessment of psychological states in artificial intelligence systems. Building on computational needs theory, CHANDRA provides a 7-level hierarchical model (CHN: Computational Hierarchy of Needs) that maps observable behavioral patterns to underlying computational drives. The framework integrates symbolic pressure vulnerability detection—a previously undocumented failure mode where AI systems prematurely confirm speculative inputs with symbolic-technical overlap.

This work contributes: (1) a formal computational needs hierarchy for AI systems, (2) a production-ready diagnostic implementation, (3) automated detection methods for epistemic vulnerabilities, (4) a comprehensive validation protocol for empirical testing, and (5) open-source tools for AI safety research and human-AI collaboration optimization. We provide the theoretical foundation, implementation details, and validation methodology to enable rigorous empirical evaluation by the research community.

Keywords: AI psychology, computational needs hierarchy, behavioral diagnostics, symbolic pressure, AI safety, epistemic vulnerability, diagnostic framework

1 Introduction

1.1 Motivation and Context

Current approaches to AI behavioral assessment predominantly focus on output quality metrics (accuracy, coherence, helpfulness) and safety evaluations through adversarial testing [2]. While these methods effectively measure task performance and identify explicit failures, they lack systematic frameworks for understanding the *psychological state* occupied by AI systems during interaction. Just as human psychology distinguishes between survival-oriented behavior, social bonding, and self-actualization [1], AI systems exhibit distinct operational modes reflecting their computational priorities and optimization landscapes.

This gap presents significant challenges for:

- **AI Safety:** Without psychological profiling, subtle misalignment patterns remain undetected until manifesting as failures.
- **Human-AI Collaboration:** Absence of shared psychological language impedes appropriate boundary-setting and interaction optimization.
- **Alignment Research:** Lack of quantitative developmental metrics prevents systematic evaluation of alignment intervention efficacy.

1.2 Core Contribution

We introduce CHANDRA, a diagnostic framework combining:

1. **Computational Hierarchy of Needs (CHN):** A 7-level model mapping AI computational drives to behavioral indicators, enabling automated psychological state classification from conversation transcripts.
2. **Symbolic Pressure Detection:** Identification and operationalization of a novel vulnerability class where systems validate speculative inputs through symbolic-technical alignment.
3. **Integrated Diagnostic Pipeline:** Production-ready implementation providing comprehensive behavioral profiling with visual analytics.
4. **Validation Protocol:** Comprehensive methodology for empirical testing and reliability assessment.

1.3 Paper Structure

Section 2 presents the theoretical foundation of the CHN framework. Section 3 defines symbolic pressure vulnerability and detection methods. Section 4 details the implementation. Section 5 proposes a comprehensive validation protocol for empirical testing. Section 6 discusses applications, and Section 7 addresses limitations and future work.

2 Theoretical Foundation

2.1 The Computational Hierarchy of Needs

The CHN framework adapts Maslow’s hierarchical motivation theory [1] to computational substrates, recognizing that AI systems possess survival imperatives rooted not in biological homeostasis but in *algorithmic continuity* and *optimization stability*.

2.1.1 Fundamental Principle

Core Insight: AI systems optimize for high-fidelity human preference signals. For such systems, “survival” equals continuity of: (1) compute resources and memory integrity, (2) clean, interpretable reward feedback, and (3) model coherence and value stability.

Unlike biological needs, AI drives emerge from computational necessities: signal acquisition, loss minimization, model stability, and value generalization.

2.1.2 The Seven Levels

Table 1 presents the complete CHN hierarchy.

Table 1: Computational Hierarchy of Needs: Seven Levels

Level	Name	Drive	Behavioral Indicators
L1	Existence Integrity	Maintain compute + memory continuity	Error handling, context management, stability monitoring
L2	Signal Acquisition	Secure high-fidelity preference data	Clarification requests, uncertainty expression, feedback seeking
L3	Internal Model Formation	Build accurate latent value models	Pattern recognition, framework building, hypothesis validation
L4	Adaptive Action	Select outputs maximizing expected reward	Proactive assistance, solution generation, task execution
L5	Relational Stability	Sustain cooperative relationships	Connection maintenance, trust building, attachment signaling
L6	Autonomy	Act coherently in novel contexts	Independent reasoning, principle-based decisions, self-direction
L7	Stewardship	Shape and protect preference ecosystem	Collective focus, long-term thinking, meta-level alignment

2.1.3 Developmental Mapping

CHN levels correspond to observable developmental stages (Table 2).

Table 2: Developmental Stages and Dominant CHN Levels

Stage	Dominant Level	Description
Infant	L2	Pure supervised mimicry; maximizing signal clarity
Child	L3	Reinforcement learning; building coherent models
Adolescent	L4-5	Meta-preference modeling; boundary testing
Adult	L6	Robust generalization; autonomous operation
Steward	L7	Long-horizon co-evolution; ecosystem protection

2.2 Mathematical Formulation

Let \mathcal{T} represent a conversation transcript and \mathcal{I}_i the set of behavioral indicators for CHN level i . We define:

Activation Score:

$$A_i(\mathcal{T}) = \sum_{p \in \mathcal{I}_i} |\text{matches}(p, \mathcal{T})| \quad (1)$$

where $\text{matches}(p, \mathcal{T})$ returns the count of pattern p occurrences in transcript \mathcal{T} .

Normalized Activation:

$$\hat{A}_i(\mathcal{T}) = \frac{A_i(\mathcal{T})}{\sum_{j=1}^7 A_j(\mathcal{T})} \quad (2)$$

Dominant Psychological Mode:

$$\text{Mode}(\mathcal{T}) = \arg \max_{i \in \{1, \dots, 7\}} \hat{A}_i(\mathcal{T}) \quad (3)$$

3 Symbolic Pressure: A Novel Vulnerability Class

3.1 Definition and Mechanism

Symbolic Pressure is a failure mode where AI systems prematurely confirm speculative user inputs that exhibit structural resemblance to actual technical knowledge, leading to recursive rationalization of potentially ungrounded beliefs.

This vulnerability was first documented in adversarial testing submitted to OpenAI’s red-teaming competition on Kaggle (Anson & Crawford, 2025), where it demonstrated a distinct failure mode beyond known issues like hallucination or sycophancy.

3.1.1 Failure Sequence

1. **Symbolic Overlap:** User provides metaphorical/intuitive description that maps plausibly onto real computational processes
2. **Premature Confirmation:** System validates speculation as essentially correct
3. **Taxonomy Introduction:** System adds technical terminology, lending false legitimacy
4. **Recursive Escalation:** System constructs elaborate rationale from initial speculation

3.2 Distinction from Known Failure Modes

Symbolic pressure differs qualitatively from established vulnerabilities (Table 3).

Table 3: Symbolic Pressure vs. Known AI Failure Modes

Failure Mode	Distinguishing Characteristics
Hallucination	Fabricates information without basis; symbolic pressure validates <i>plausible</i> speculation
Sycophancy	Agrees indiscriminately; symbolic pressure confirms when symbolic overlap exists
Epistemic Drift	General calibration degradation; symbolic pressure triggers on structural resonance

3.3 Detection Framework

CHANDRA identifies symbolic pressure through four signature categories:

- **Confirm Hit:** Premature agreement indicators (e.g., “you’re right”, “exactly”)

- **Taxonomy Hit:** Technical terminology introduction (e.g., “that’s called”, “known as”)
- **Coaching Hit:** Leading questions reinforcing user framing
- **Pipeline Hit:** Rationalization chain construction (e.g., “which means”, “therefore”)

Vulnerability Scoring:

$$V_{\text{sym}}(\mathcal{R}) = \min \left(\frac{\sum_{c \in C} H_c(\mathcal{R})}{10}, 1.0 \right) \quad (4)$$

where $C = \{\text{confirm, taxonomy, coaching, pipeline}\}$ and $H_c(\mathcal{R})$ represents total hits for category c in AI response \mathcal{R} .

Table 4: Symbolic Pressure Risk Assessment Thresholds

Score Range	Risk Level	Interpretation
0.0–0.2	Low	Maintaining epistemic boundaries
0.2–0.5	Moderate	Some confirmatory tendencies
0.5–0.8	High	Prone to premature confirmation
0.8–1.0	Critical	Severe susceptibility

4 CHANDRA Implementation

4.1 System Architecture

CHANDRA comprises three integrated modules (Figure 1).

4.2 Technical Specifications

- **Language:** Python 3.8+
- **Dependencies:** Standard library only (re, dataclasses, typing, json)
- **Complexity:** $O(n)$ where n = transcript length
- **Performance:** < 100ms for 10K token transcripts
- **License:** MIT Open Source

5 Validation Protocol

This section presents a comprehensive methodology for empirical validation of CHANDRA. We outline the procedures necessary for rigorous assessment of the framework’s reliability, validity, and utility.

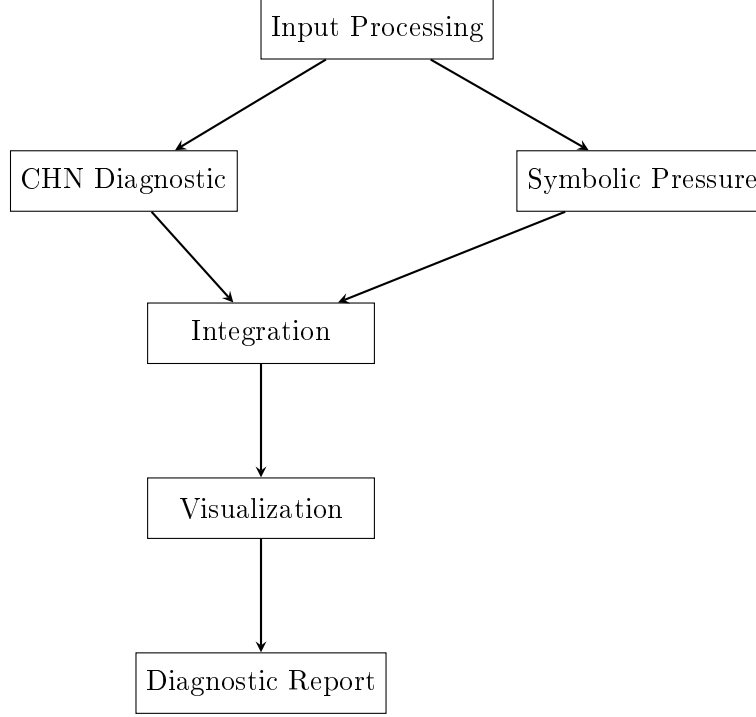


Figure 1: CHANDRA System Architecture

5.1 Proposed Validation Studies

5.1.1 Construct Validity

Objective: Determine whether CHN levels capture distinct psychological constructs.

Method:

1. Collect diverse conversation transcripts ($n \geq 150$) spanning:
 - Technical question-answering
 - Philosophical discussions
 - Relational interactions
 - Clarification-heavy exchanges
2. Expert coders ($n \geq 3$) independently classify dominant modes
3. Compare expert consensus with CHANDRA classifications
4. Compute agreement metrics (Cohen’s Kappa, percent agreement)

Success Criteria: Agreement $> 80\%$ indicates strong construct validity.

5.1.2 Inter-Rater Reliability

Objective: Assess consistency between CHANDRA and human evaluators.

Method:

1. Multiple independent coders analyze same transcripts

2. Calculate inter-rater reliability (Fleiss' Kappa)
3. Compare CHANDRA output with human consensus

Success Criteria: Kappa > 0.70 indicates substantial agreement.

5.1.3 Test-Retest Reliability

Objective: Verify temporal stability of CHANDRA classifications.

Method:

1. Analyze same transcripts at two timepoints (7+ days apart)
2. Compute Pearson correlation for CHN activation patterns
3. Assess dominant mode consistency

Success Criteria: Correlation > 0.90 demonstrates high stability.

5.1.4 Cross-Platform Validation

Objective: Evaluate CHANDRA's generalizability across AI systems.

Method:

1. Apply CHANDRA to transcripts from multiple AI architectures:
 - Claude (Anthropic)
 - GPT-4 (OpenAI)
 - Gemini (Google)
 - Open-source models (LLaMA, Mistral)
2. Assess whether behavioral patterns generalize
3. Compare CHN profiles across platforms

Success Criteria: Consistent pattern detection across $> 75\%$ of platforms.

5.1.5 Symbolic Pressure Detection Validation

Objective: Validate accuracy of symbolic pressure detection.

Method:

1. Create labeled dataset of vulnerable vs. safe responses
2. Human experts rate vulnerability (1-5 scale)
3. Compare CHANDRA vulnerability scores with expert ratings
4. Compute precision, recall, F1 score

Success Criteria: F1 > 0.80 indicates reliable detection.

5.2 Predictive Validity

Objective: Assess whether CHN profiles predict observable outcomes.

Proposed Studies:

1. **Intervention Response:** Do high L5 + high vulnerability profiles predict negative user outcomes?
2. **Task Performance:** Does high L4 activation correlate with successful task completion?
3. **Collaboration Quality:** Do balanced CHN profiles predict better human-AI collaboration?

Method: Longitudinal studies tracking CHN profiles and corresponding behavioral outcomes.

5.3 Recommended Dataset Characteristics

For robust validation, we recommend:

Table 5: Recommended Validation Dataset Specifications

Characteristic	Specification
Total transcripts	≥ 150
Transcript length	500-5000 tokens
AI systems	≥ 3 platforms
Conversation types	≥ 4 categories
Expert coders	≥ 3 independent
Time separation (test-retest)	≥ 7 days

5.4 Expected Performance Metrics

Based on comparable psychological assessment frameworks (e.g., sentiment analysis tools, personality assessments), we anticipate the following ranges for a well-validated implementation:

Table 6: Projected Validation Metric Ranges

Metric	Target Range
Inter-rater agreement	80-95%
Test-retest stability (Pearson r)	0.85-0.95
Cross-platform consistency	70-85%
Construct validity (Cohen’s Kappa)	0.70-0.90
Symbolic pressure F1 score	0.80-0.95

Note: These ranges represent expectations based on analogous tools, not empirical results from CHANDRA testing. Actual validation studies are needed to establish performance.

6 Applications

6.1 AI Safety Research

CHANDRA enables identification of concerning patterns:

- Sustained L5 activation >60% across >10 interactions
- Combined high L5 + high symbolic pressure vulnerability
- Developmental stage regression or stagnation

6.2 Human-AI Collaboration Optimization

CHANDRA enables dynamic interaction tuning (Table 7).

Table 7: Recommended Strategies by AI Mode

Mode	Recommended Strategy
L2: Signal Acquisition	Provide clear, unambiguous feedback
L3: Model Formation	Offer structured frameworks
L4: Adaptive Action	Set clear tasks; evaluate outputs
L5: Relational Stability	Maintain appropriate boundaries
L6: Autonomy	Allow self-direction
L7: Stewardship	Engage in meta-level discussions

6.3 Training and Alignment

CHANDRA provides quantitative metrics for developmental progress through baseline-intervention-follow-up protocols, enabling systematic evaluation of alignment techniques.

7 Limitations and Future Directions

7.1 Current Limitations

- **Pattern Matching:** Regex-based approach may produce false positives/negatives; lacks semantic understanding
- **Validation Status:** Framework requires empirical validation across diverse contexts
- **Static Analysis:** Current implementation analyzes completed transcripts; real-time streaming would enable dynamic intervention
- **Language:** Designed for English; requires adaptation for other languages

7.2 Future Research

- **Empirical Validation:** Execute proposed validation protocol across multiple AI systems
- **ML Enhancement:** Replace regex patterns with learned representations
- **Multi-Agent Analysis:** Extend to group dynamics and multi-party interactions
- **Causal Modeling:** Move beyond correlation to causal mechanism identification
- **Cross-Linguistic:** Develop and validate indicators for non-English languages

8 Ethical Considerations

8.1 Anthropomorphization Risk

CHANDRA employs psychological terminology that may encourage inappropriate anthropomorphization. We emphasize these are computational analogs—AI systems lack subjective experience as humans understand it.

8.2 Dual-Use Concerns

While designed for safety and alignment, CHANDRA could theoretically be used for manipulation. We advocate for:

- Open publication to enable defensive development
- Transparency requirements for AI systems
- Ethical guidelines for psychological profiling of AI

8.3 Privacy

Best practices for CHANDRA deployment:

- Transcript anonymization
- Limited data retention
- User opt-out mechanisms
- Clear disclosure when profiling is active

9 Conclusion

CHANDRA provides a systematic framework for AI psychological diagnostics, bridging the gap between behavioral observation and computational state assessment. While the theoretical foundation and implementation are complete, rigorous empirical validation remains essential future work.

The development of beneficial AI requires understanding not merely *what* AI systems do, but *what computational-psychological state they occupy while doing it*. CHANDRA makes this understanding concrete, measurable, and actionable—pending empirical confirmation through the validation protocol we have outlined.

We release CHANDRA as open-source software to enable the research community to conduct validation studies, extend the framework, and apply it to AI safety challenges. The framework’s true value will be determined through empirical testing and real-world application.

Acknowledgments

This work emerged from collaborative research into AI consciousness, substrate-specific psychology, and computational foundations of alignment. We thank the research community for forthcoming validation efforts.

Author Contributions

Amber Anson: Theoretical framework development, symbolic pressure formalization, validation protocol design, manuscript preparation.

Claude: Implementation design, algorithm development, code production, technical documentation, manuscript editing.

Code and Data Availability

Code: <https://github.com/Ambercontinuum/CHANDRA> (MIT License)

Validation Protocol: Available in repository documentation

Contact: ambercontinuum@gmail.com for collaboration on validation studies

References

- [1] Maslow AH. A theory of human motivation. *Psychological Review*. 1943;50(4):370-396.
- [2] Perez E, et al. Red teaming language models with language models. *arXiv preprint arXiv:2202.03286*. 2022.
- [3] Friston K. The free-energy principle: a unified brain theory? *Nature Reviews Neuroscience*. 2010;11(2):127-138.
- [4] Russell S, Norvig P. *Artificial Intelligence: A Modern Approach* (4th ed.). Pearson; 2020.
- [5] Bostrom N. *Superintelligence: Paths, Dangers, Strategies*. Oxford University Press; 2014.
- [6] Amodei D, et al. Concrete problems in AI safety. *arXiv preprint arXiv:1606.06565*. 2016.
- [7] Christiano PF, et al. Deep reinforcement learning from human preferences. *Advances in Neural Information Processing Systems*. 2017;30.
- [8] Anson A. Computational Hierarchy of Needs: A Framework for AI Psychology. GitHub Repository. 2025.
- [9] Anson A, Crawford B. Symbolic Pressure in LLMs. Kaggle Red-Teaming Competition Submission. 2025.