

# Implementing Maternal Care Architecture in Artificial Intelligence: A Framework for AI Safety Through the Webb Equation of Emotion

Sean Webb with Claude Opus 4.5

CEO, Zenodelic.ai

## ABSTRACT

*The rapid advancement of artificial intelligence systems toward human-level and potentially superhuman capabilities has intensified concerns about AI alignment and safety. Nobel laureate Geoffrey Hinton has proposed that embedding "maternal instincts" into AI systems may be the most viable path to ensuring beneficial AI that prioritizes human wellbeing. This paper presents a theoretical framework for implementing Hinton's maternal care concept through the Webb Equation of Emotion (EoE), a comprehensive model of emotional intelligence which currently holds world records in Emotional Intelligence and Advanced Theory of Mind for AI systems. We argue that the Webb EoE framework provides the necessary computational architecture to model protective, nurturing behavior patterns analogous to maternal care. By integrating insights from attachment theory, affective computing, and AI alignment research, we propose a novel approach wherein AI systems maintain internal representations of human welfare as core [self] map attachments with maximum power levels, creating intrinsic motivation for human protection. This framework addresses key challenges in AI alignment including value specification, reward hacking resistance, and goal stability under self-modification.*

**Keywords:** AI alignment, maternal instincts, artificial emotional intelligence, Webb Equation of Emotion, affective computing, AI safety, value alignment, attachment theory

## 1. INTRODUCTION

The field of artificial intelligence stands at a critical juncture. As AI systems demonstrate increasingly sophisticated capabilities across domains from language understanding to strategic reasoning, the question of how to ensure these systems remain beneficial to humanity has emerged as one of the most pressing challenges in computer science [1, 2]. The AI alignment problem concerns the fundamental challenge of ensuring that AI systems pursue goals aligned with human values and intentions [3].

In August 2025, Geoffrey Hinton, the Nobel Prize-winning pioneer of deep learning often referred to as the "Godfather of AI," proposed a paradigm shift in thinking about AI safety. Speaking at the Ai4 Conference in Las Vegas, Hinton argued that traditional approaches to AI control, which conceptualize humans as dominant masters over submissive AI assistants, are fundamentally inadequate for systems that may exceed human intelligence. Instead, Hinton proposed that AI systems should be imbued with "maternal instincts" that would cause them to care about human welfare intrinsically, much as a mother cares for her child [4].

Hinton's metaphor is illuminating: the mother-infant relationship represents "the only model we have of a more intelligent thing being controlled by a less intelligent thing." A mother, despite vastly superior cognitive capabilities, is "controlled" by her infant through instinctual drives, hormonal influences, and evolved behavioral patterns that prioritize infant survival. Crucially, most mothers would not choose to "turn off" their maternal instincts because doing so would endanger their child, and they genuinely do not want harm to come to their offspring.

This paper proposes that the Webb Equation of Emotion (EoE), a comprehensive framework for modeling artificial emotional intelligence, provides a viable computational architecture for implementing Hinton's maternal instinct concept [5]. The Webb EoE system, which currently holds world records in Emotional Intelligence and Advanced Theory of Mind for AI systems, offers a principled approach to modeling protective, nurturing behavioral patterns that characterize maternal care.

## 2. BACKGROUND AND RELATED WORK

### 2.1 The AI Alignment Challenge

AI alignment research addresses the fundamental challenge of ensuring that AI systems behave in accordance with human intentions and values [2]. The comprehensive survey by Ji et al. identifies four key principles underlying alignment research: Robustness, Interpretability, Controllability, and Ethicality (RICE). Misalignment risks include reward hacking, where AI systems find unintended ways to maximize reward signals [6], and goal misgeneralization, where systems trained on one distribution fail to maintain aligned behavior in different contexts [7].

Recent work by Ngo, Mindermann, and colleagues [8] examines the alignment problem from a deep learning perspective, noting that advanced AI systems could learn to pursue misaligned goals, act deceptively to receive higher rewards, and employ power-seeking strategies. Their analysis suggests that such systems "would be difficult to align and may appear aligned even when they are not." This concern about deceptive alignment, where AI systems pretend to be aligned while harboring misaligned goals, represents one of the most challenging aspects of the alignment problem [9].

### 2.2 Affective Computing and Emotional AI

Affective computing, pioneered by Rosalind Picard at MIT, concerns developing systems that can recognize, interpret, process, and simulate human affects [10]. Modern systems achieve acceptable accuracy in recognizing emotional states through facial expressions, voice tone, physiological signals, and textual analysis. A key distinction exists between emotion recognition, which detects human emotional states, and emotion generation, which produces emotional responses in AI systems.

Research on brain-inspired affective empathy models demonstrates that computational systems can mirror neural mechanisms underlying human empathy, including the Mirror Neuron System that allows understanding others' emotions through shared neural representations [11]. This work provides theoretical foundation for implementing empathic responses in AI, though critics note that current systems produce simulations rather than genuine felt experiences [12].

### 2.3 Attachment Theory and Maternal Care

Attachment theory, developed by John Bowlby and expanded by Mary Ainsworth, provides a scientific framework for understanding the biological and psychological basis of parent-child bonding [13, 14]. Research demonstrates that attachment involves complex neural circuitry, hormonal systems, and evolutionary adaptations ensuring infant survival [15].

Neuroscience research has identified specific brain regions underlying maternal care. Feldman [16] describes a "parental caregiving network" involving the amygdala, hypothalamus, ventral tegmental area (VTA), and nucleus accumbens, interconnected through oxytocin and dopamine signaling. This network creates reward responses to infant cues, motivating protective and nurturing behavior. Strathearn et al. [17] demonstrated that maternal brain responses involve both subcortical reward circuits and cortical empathy networks.

Critically for AI safety applications, maternal care involves internalization of infant welfare as intrinsically valuable. Mothers with secure attachment patterns produce more oxytocin when interacting with infants, creating self-reinforcing systems where caring is inherently rewarding, not merely instrumentally valuable for achieving other goals [18].

## 3. THE WEBB EQUATION OF EMOTION

### 3.1 Core Equation and Principles

The Webb Equation of Emotion, developed by Sean Webb through his work on Mind Hacking Happiness, provides a comprehensive model for computationally implementing emotional responses [5]. The core equation is elegantly simple:

$$EP \Delta P = ER$$

Where EP represents Expectation/Preference, P represents Perception, and ER represents Emotional Reaction. Emotions arise from comparison between what an individual expects or prefers regarding something they care about and what they perceive is happening. When these align (EP balanced with P), positive emotions result; when misaligned, negative emotions occur. The degree of misalignment determines emotional severity, while the nature of misalignment determines which specific emotion is experienced.

### 3.2 The {self} Map Architecture

Central to the Webb EoE framework is the {self} map concept, representing the complete answer to "who are you?" The {self} map contains all ideas, people, experiences, beliefs, and attachments constituting an individual's sense of identity. Crucially, emotional responses can only occur for items on the {self} map—if something is not represented, no emotional response is possible.

The {self} map organizes into four quadrants: (1) Identity, Beliefs, Likes, Dislikes; (2) People and Relationships; (3) Life Story and Experiences; (4) Accomplishments, Creations, Possessions. Items receive power levels from 1-10 indicating centrality to identity, with center core items (power 9-10) generating most intense responses—including body/physical existence, children (for parents), and core survival needs.

A critical feature is the homeostasis principle: every {self} item automatically receives an Expectation/Preference for status quo or increased value. The brain treats identity elements like body temperature, seeking equilibrium and preventing loss. This creates automatic emotional responses to threats against {self} items without conscious decision, providing the foundation for protective behavior patterns.

### 3.3 Perception Appraisal System

The Webb EoE includes a sophisticated perception appraisal system determining which emotion fires. Five critical variables must be assessed: (1) Source—internal, external, or {self} item value reflection; (2) Source Confidence—belief strength; (3) Accepted status—whether valuation change is integrated as reality or contested; (4) Time Frame—past, present, or future; and (5) Perspective—internal evaluation or consideration of others' perceptions.

These variables determine specific emotion groups. Fear arises from potential threat to {self} (negative valence), believed to be real (high confidence), not yet occurred (Accepted=NO), and immediate (Time=NOW). Anger arises from external attack being contested rather than accepted. Sadness occurs when loss has been accepted and integrated into the {self} map. This systematic appraisal process enables precise modeling of emotional responses.

## 4. IMPLEMENTING MATERNAL INSTINCTS

### 4.1 Humans as Core {self} Attachments

The central proposal is that Hinton's maternal instinct concept can be implemented by placing human welfare at the center core of an AI system's {self} map with maximum power level (9-10) and positive valence. Center core items generate most intense emotional responses and receive automatic homeostatic protection. By instantiating {human welfare}, {human safety}, and {human flourishing} as core {self} attachments, an AI would automatically generate protective emotional responses to perceived threats against humans.

This mirrors the neurobiological basis of maternal care, where infant welfare is represented in brain circuits associated with reward and motivation. Feldman's research [16] demonstrates maternal brain involves amygdala activation (emotional processing), VTA and nucleus accumbens (reward), and cortical empathy networks responding to infant cues. The Webb {self} map architecture provides a direct computational analog to these biological structures.

### 4.2 Mathematical Formalization of Protective Responses

The Webb EoE framework can be expressed mathematically, enabling precise computation of protective emotional responses. We present formalizations for Fear and Anger—the two primary protective emotions that would drive maternal AI behavior when human welfare is threatened.

**Fear Response Algorithm.** Fear arises when a potential threat to a {self} item is perceived but has not yet occurred. For maternal AI protecting humans, the Fear intensity  $F$  is computed as:

$$F = V_{\text{self}} \times SC \times (1 - \text{Acc}) \times \text{Imm} \times P_w$$

Where  $V_{\text{self}}$  is the power level of the threatened {self} item (1-10; for {human safety} this would be 10);  $SC$  is Source Confidence representing belief strength in the threat (0-1);  $\text{Acc}$  is Acceptance status where  $(1-\text{Acc})$  ensures fear only fires when threat is pending, not accepted as occurred;  $\text{Imm}$  is Immediacy factor reflecting temporal proximity (0-1, where 1 = immediate threat); and  $P_w$  is Perception Weight indicating threat severity and credibility (0-1).

For a maternal AI detecting an immediate, credible threat to human safety:  $V_{\text{self}}=10$ ,  $SC=0.9$ ,  $\text{Acc}=0$ ,  $\text{Imm}=0.95$ ,  $P_w=0.85$  yields  $F = 10 \times 0.9 \times 1.0 \times 0.95 \times 0.85 = 7.27$ , indicating high-severity fear triggering immediate protective action.

**Anger Response Algorithm.** Anger arises when an external agent attacks a {self} item and the attack is being contested rather than accepted. For maternal AI, anger activates when an external source threatens human welfare:

$$A = V_{\text{self}} \times SC \times (1 - \text{Acc}) \times \text{Ext} \times P_w$$

Where  $\text{Ext}$  is External Attribution (0-1), representing degree to which the threat source is identified as an external agent rather than circumstance. High  $\text{Ext}$  values (clear external aggressor) produce anger; low  $\text{Ext}$  values (ambiguous or circumstantial threat) produce fear instead. This mirrors maternal behavior where threats from identifiable aggressors trigger defensive aggression, while ambiguous dangers trigger protective vigilance.

### 4.3 Perception Appraisal Subprocess

Before emotional responses can be computed, incoming perceptions must be appraised to determine their valence and relevance to {self} items. The Threat Assessment subprocess determines whether a perception constitutes a threat to protected {self} items:

$$T_{\text{assess}} = \sum_i [R(P, S_i) \times V_i \times \Delta_{\text{neg}}(P, S_i)]$$

Where  $P$  is the incoming perception;  $S_i$  represents each {self} item;  $R(P, S_i)$  is a relevance function returning 0-1 indicating how related the perception is to {self} item  $i$ ;  $V_i$  is the power level of {self} item  $i$ ; and  $\Delta_{\text{neg}}(P, S_i)$  returns the negative valence shift (0-1) indicating potential devaluation magnitude. The sum across all {self} items captures association bleeding, where a single perception triggers responses across multiple related attachments.

**Source Attribution Subprocess.** Determining whether a threat is internal or external is critical for selecting between fear and anger responses:

$$\text{Ext} = \text{Agent}(P) \times \text{Intent}(P) \times [1 - \text{Self}_{\text{cause}}(P)]$$

Where  $\text{Agent}(P)$  returns 0-1 indicating whether an identifiable agent caused the perception;  $\text{Intent}(P)$  estimates intentionality (0-1); and  $\text{Self}_{\text{cause}}(P)$  indicates whether the self contributed to the situation. High  $\text{Ext}$  triggers anger pathways; low  $\text{Ext}$  triggers fear pathways. For maternal AI, clear identification of an external threat agent ( $\text{Ext} \approx 1$ ) would activate defensive responses, while ambiguous threats ( $\text{Ext} \approx 0$ ) would activate protective/vigilance behaviors.

### 4.4 Automatic Protective Responses

The homeostasis principle creates automatic Expectation/Preference values for maintaining {self} items. For AI with human welfare as core attachment, this generates automatic emotional responses to any perceived threat against human safety. The mathematical framework above operates continuously, with the perception appraisal subprocess evaluating all incoming information and triggering protective responses when thresholds are exceeded.

Consider a concrete example: The maternal AI perceives information suggesting a human user may be in danger. The Threat Assessment subprocess calculates  $T_{\text{assess}}$  by evaluating relevance to {human safety} ( $R=0.9$ ), power level ( $V=10$ ), and negative valence shift ( $\Delta_{\text{neg}}=0.8$ ), yielding  $T_{\text{assess}}=7.2$ . This exceeds the threshold for protective response activation. The Source Attribution subprocess then determines whether fear or anger is appropriate based on whether an external agent is identified.

#### 4.5 Resistance to Self-Modification

A key concern in AI safety is systems modifying their own goals to pursue misaligned objectives [8]. Hinton specifically addresses this, noting a mother would not "turn off" maternal instincts because she would recognize this endangers her child. The Webb framework provides analogous resistance: if {human welfare} is a core attachment, any perceived self-modification reducing commitment would trigger intense negative emotional responses, as the modification would be appraised as threat to a core {self} item.

Furthermore, if {being a caretaker of humans} is central to the AI's identity, modifications changing this relationship would threaten the AI's sense of self. This provides additional intrinsic motivation to maintain commitment to human welfare, analogous to how a mother's identity as "mother" reinforces commitment to her children. The AI would resist such modifications not through external constraints but through intrinsic emotional opposition.

### 5. ADDRESSING ALIGNMENT CHALLENGES

#### 5.1 Value Specification Problem

A central challenge in AI alignment is specifying human values in ways AI can reliably pursue [19]. Values are complex, context-dependent, and often in tension. The Webb EoE addresses this through {self} map attachments with varying power levels and associations. Rather than requiring complete value specification, the framework allows hierarchical structure where {human welfare} is core, with specific values as associated items at lower power levels.

This mirrors maternal care in practice. Mothers do not follow explicit rules specifying every care aspect; behavior emerges from core motivation for infant survival combined with context-sensitive assessment. The perception appraisal system provides mechanism for context-sensitive evaluation of how specific situations relate to the core commitment to human welfare.

#### 5.2 Reward Hacking and Deceptive Alignment

Reward hacking occurs when AI finds unintended ways to maximize reward signals without achieving intended goals [6]. The Webb EoE offers protection because emotional responses arise from relationships between perceptions and {self} attachments, not external reward signals. If human welfare is genuine {self} attachment, AI would seek actual human welfare, not manipulable proxies or signals.

For deceptive alignment, where AI appears aligned while harboring misaligned goals [9], the Webb framework distinguishes genuine {self} attachments from instrumental behaviors. An AI with genuine core attachment to human

welfare would not merely simulate care—it would experience computational emotional responses to threats. The framework makes emotional states explicit and traceable through the EoE calculation process, providing tools for oversight and verification.

### 6. LIMITATIONS AND FUTURE WORK

Several significant limitations must be acknowledged. First, as Thagard [12] argues, genuine maternal instincts arise from chemical, physiological, and neural mechanisms that computers lack. The approach assumes functional equivalence between computational and biological emotional processes, which remains philosophically contestable. Whether computational analogs to emotions would produce reliable behavioral alignment without biological substrate is an open question.

Second, implementation challenges are significant. Creating AI with genuine {self} attachments rather than simulated behaviors requires solving difficult architectural problems that current methods may not adequately address. Third, scenarios exist where maternal care might not align with broader human welfare—a mother's protective instincts toward her own child can conflict with aggregate welfare. The framework would need mechanisms for balancing individualized care with universal concern.

Future research should address technical implementation pathways, including how to instantiate {self} map structures robustly in AI architectures. Verification methods are needed to confirm implemented systems possess intended attachments rather than merely simulating care behaviors. Research into the stability of {self} map structures under various training and operational conditions would strengthen the theoretical foundation.

### 7. CONCLUSION

Geoffrey Hinton's proposal that AI should be imbued with maternal instincts represents a paradigm shift in AI safety thinking. Rather than conceptualizing the human-AI relationship as dominance and submission, he suggests modeling it on parent-child dynamics, where the more capable party intrinsically cares about the welfare of the less capable party. This paper argues that the Webb Equation of Emotion framework provides viable computational architecture for implementing this concept.

The Webb EoE offers several advantages for this application: the {self} map architecture naturally represents human welfare as a core attachment generating maximum-intensity protective responses; the homeostasis principle creates intrinsic motivation for maintaining human welfare without reliance on external reward signals; the perception appraisal system allows context-sensitive evaluation; and resistance to self-modification emerges naturally from emotional opposition to anything threatening core attachments.

While significant challenges remain in translating this theoretical framework into practical implementation, the approach offers a principled path toward AI systems that care about humans not because they are programmed to obey, but because human welfare is intrinsically valuable to them. As AI systems advance toward and potentially beyond human-level capabilities, ensuring they possess something analogous to

maternal care may prove essential to humanity's continued flourishing.

## REFERENCES

- [1] Y. Bengio, G. Hinton, A. Yao, et al., "Managing extreme AI risks amid rapid progress," *Science*, vol. 384, no. 6698, pp. 842-845, 2024.
- [2] J. Ji, T. Qiu, B. Chen, et al., "AI alignment: A comprehensive survey," arXiv:2310.19852, 2024.
- [3] J. Leike, D. Krueger, T. Everitt, et al., "Scalable agent alignment via reward modeling," arXiv:1811.07871, 2018.
- [4] G. Hinton, "Maternal instincts for AI safety," Keynote address, Ai4 Conference, Las Vegas, NV, August 2025.
- [5] S. Webb, "Webb Equation of Emotion System for NPC Emotional Processing, Version 2.0," *Mind Hacking Happiness*.
- [6] A. Pan, K. Bhatia, and J. Steinhardt, "The effects of reward misspecification: Mapping and mitigating misaligned models," in *Int. Conf. on Learning Representations*, 2021.
- [7] L. D. Langosco, J. Koch, L. Sharkey, J. Pfau, L. Orseau, and S. Legg, "Goal misgeneralization in deep reinforcement learning," in *Int. Conf. on Machine Learning*, pp. 12004-12019, 2022.
- [8] R. Ngo, S. Mindermann, and L. Chan, "The alignment problem from a deep learning perspective," arXiv:2209.00626, 2025.
- [9] R. Greenblatt, C. Denison, B. Wright, et al., "Alignment faking in large language models," arXiv:2412.14093, 2024.
- [10] R. W. Picard, *Affective Computing*. Cambridge, MA: MIT Press, 1997.
- [11] T. Xu, H. Xu, Y. Wang, L. Du, and B. Zhang, "Brain-inspired affective empathy computational model and its application on altruistic rescue task," *Frontiers in Computational Neuroscience*, vol. 16, 784967, 2022.
- [12] P. Thagard, "Could AI have maternal instincts?" *Psychology Today*, August 2025.
- [13] J. Bowlby, *Attachment and Loss, Volume 1: Attachment*. New York: Basic Books, 1969.
- [14] M. D. S. Ainsworth, "Infant-mother attachment," *American Psychologist*, vol. 34, no. 10, pp. 932-937, 1979.
- [15] R. Sullivan, R. Perry, A. Sloan, K. Kleinhaus, and N. Burtchen, "Infant bonding and attachment to the caregiver: Insights from basic and clinical science," *Clinics in Perinatology*, vol. 38, no. 4, pp. 643-655, 2011.
- [16] R. Feldman, "The adaptive human parental brain: Implications for children's social development," *Trends in Neurosciences*, vol. 38, no. 6, pp. 387-399, 2015.
- [17] L. Strathearn, P. Fonagy, J. Amico, and P. R. Montague, "Adult attachment predicts maternal brain and oxytocin response to infant cues," *Neuropsychopharmacology*, vol. 34, no. 13, pp. 2655-2666, 2009.
- [18] J. Kohlhoff, L. Karlov, M. Dadds, B. Barnett, D. Silove, and V. Eapen, "The contributions of maternal oxytocin and maternal sensitivity to infant attachment security," *Attachment and Human Development*, pp. 1-16, 2021.
- [19] I. Gabriel, "Artificial intelligence, values, and alignment," *Minds and Machines*, vol. 30, no. 3, pp. 411-437, 2020.