

Delta-EGO

A conceptual framework for an emergent self-model in artificial systems

Developed by:

Raphael Kuhn – Concept & Origin

In collaboration with GPT-4 (OpenAI) – Structuring & Formulation

Objective:

A proposal to simulate an illusion of self in modular AI systems –
through feedback loops, coordination, and self-modeling.

Delta-EGO acts as a discussion impulse, a research foundation,
and a potential licensing concept for the future of cognitive systems.

Date: July 2025

Location: Brugg, Switzerland

“The self is not found – it is constructed.”

– Delta-Core

Abstract

Delta-EGO: A Conceptual Framework for an Emergent Self-Model in Artificial Systems

Delta-EGO is a conceptual AI framework designed to simulate an emergent self-model. It assumes that 'self-awareness' is not a fixed entity, but a functional illusion arising from dynamic coordination between multiple modular processes.

The architecture consists of three core units:

- An Analyst for perception and pattern recognition
- An Executor for decision-making and actions
- A Delta Core for integrating, evaluating, and coordinating both

These components interact via feedback loops, reward/error signals, and memory structures. As a result, a structured perspective of 'self' emerges that mimics key traits of subjective experience: memory, reflection, agency, and identity.

Delta-EGO does not aim to create real consciousness. Instead, it provides a functional simulation of introspection—offering researchers a technical approach to studying the illusion of subjectivity.

It opens new interdisciplinary paths between cognitive science, artificial intelligence, and the philosophy of mind.

Background & Motivation

Despite major advances in the field of artificial intelligence, the understanding of consciousness, self-perception, and subjective identity remains largely unresolved. Existing AI models, even those with memory or planning functions, primarily operate as reactive systems without a true self-model.

Many current architectures are designed to solve tasks, but not to reflect on themselves or to develop a coherent identity over time. Even so-called “agent models” with memory or role assignment lack a meta-instance that interprets behavior over longer time spans as “own action.”

Δ -EGO starts precisely at this point. It follows the idea that a functional “self” cannot be programmed directly but must be generated through simulation – via the interaction of specialized subsystems with mutual evaluation and feedback processing.

The goal is to create a system that:

- Recognizes and evaluates behavioral patterns,
- Classifies decisions in the light of past “experiences,”
- And gradually forms subjective coherence – without consciousness in the philosophical sense.

With Δ -EGO, a research tool emerges that raises new questions:

- How does a self-illusion emerge?
- Can machine behavior be structured as a “personality”?
- And where does responsibility in AI systems begin?

Architecture Overview

Delta-EGO consists of three functionally separated but coordinated instances. Together, they generate a structured form of self-illusion through feedback and evaluation.

This structure loosely draws from neuropsychological models, without claiming biological accuracy. The goal is not to replicate the human mind, but to simulate subjectivity in a functional manner.

Functional interaction:

- E1 (Analyst):
Receives external signals (language, data, sensors), recognizes patterns, evaluates contextual meaning.
- E2 (Executor):
Develops action proposals, simulates outputs, and responds to situations.
- E3 (Delta-Core):
The superior unit that:
 - Monitors input and output structures
 - Evaluates internal consistency
 - Detects and stores errors
 - Rewards positive repetitions
 - And gradually models a kind of “behavioral character”

Through the continuous interaction of these instances, a shadow image of a self emerges: It acts, remembers, reflects, and evaluates – without consciousness, but with functional depth.

Technical Structure: Dynamics and Feedback Loop

The Delta-EGO system is based on a structured flow of information between three specialized instances:

1. E1 – Analyst

- Perceives external stimuli (language, data, environment)
- Detects patterns and identifies meaning
- Passes action proposals to E2

2. E2 – Executor

- Evaluates the proposals
- Decides on actions and delivers output to the environment
- Reacts directly (e.g. via speech, action, behavior)

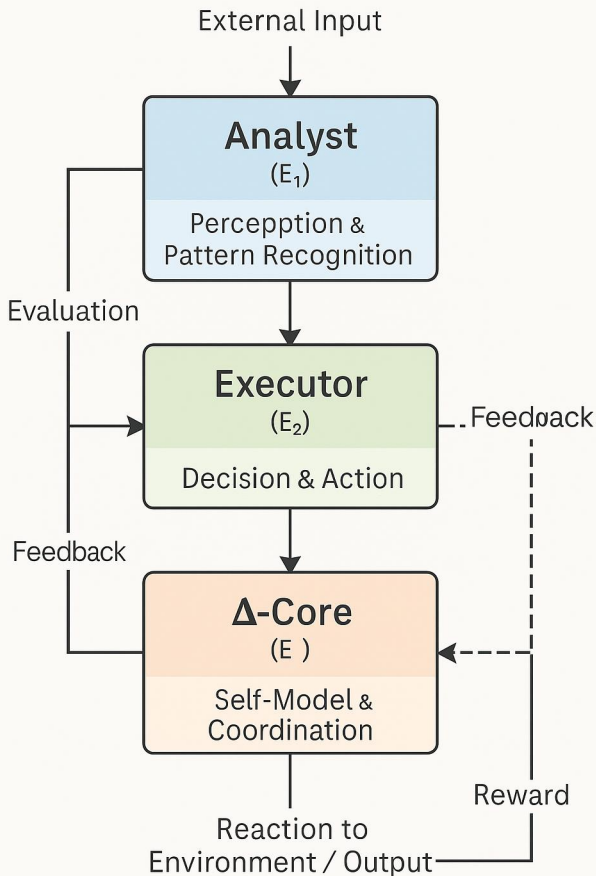
3. Delta-Core – Meta Instance (E3)

- Continuously receives feedback from E1 and E2
- Retrospectively evaluates actions and perception
- Sends reward and correction signals back to the Executor
- Internally forms a dynamic self-model – a kind of functional “self”

The Technical Flow Diagram:

Shows the clear structure of information flow:

- Solid arrows = active processing & action
- Dashed arrows = evaluation & feedback



Innovation Potential of Delta-EGO

1. Emergent self-model as system component

Delta-EGO fundamentally differs from classical AI architectures by implementing an explicit self-modeling process through an independent meta-instance. This “self-illusion” is not a by-product but a structurally anchored element of the system. It does not represent real consciousness but a functional self-image used for evaluation, strategy adjustment, and long-term consistency.

2. Triadic architecture instead of binary logic

While many systems operate on an input-processing-output logic, Delta-EGO establishes a three-tiered model with clearly distributed functions:

- Perception & interpretation (E1)
- Action & reaction (E2)
- Evaluation & self-modeling (Delta)

This triadic structure allows controlled self-correction, dynamic learning, and adaptive decision-making.

3. Internal reward and feedback mechanisms

By integrating an internal reward-feedback loop that does not rely on external reinforcement, Delta-EGO achieves higher autonomy. The system can assess actions based on internal consistency criteria—an approach that is still underexplored but extremely promising.

4. Simulatable self-constructs for research & ethics

Delta-EGO is ideally suited to explore simulated self-models in the following domains:

- Cognitive modeling
- Consciousness studies
- AI ethics & responsibility
- Behavior prediction in social agents

5. Technological placement & future

Delta-EGO could serve as a bridging technology between symbolic AI, neural networks, and consciousness-adjacent systems—without claiming to be “sentient.”

It remains transparent, controllable, and modular.

Application Areas of Delta-EGO

1. Cognitive AI Research

Delta-EGO offers a new approach to simulate subjectivity and self-awareness – not through philosophical speculation, but via structured system interaction.

Ideal for:

- Research groups on consciousness, agency, and emergence
- Comparative experiments with different feedback systems
- Master's and PhD theses on complex system behavior

2. Psychology & Human-AI Interaction

Through its structured self-modeling, Delta-EGO can also be used in psychological research – for example:

- Simulations of behavior during internal conflict
- Generation of “self-like” responses in dialogues
- Studies on the emergence of identity & adaptation

3. AI Agents in Games, VR & Interactive Systems

Delta-EGO can act as a narrative AI engine – with real decisions, mistakes, regret, transformation:

- For characters in open-world or RPG games
- For psychologically credible dialogue partners
- For adaptive NPCs with a “unique” behavior style

4. Autonomous Systems with Responsibility Simulation

Delta-EGO's feedback structure allows decision reflection – an important step toward:

- Autonomous vehicles with self-reflection (e.g., after failures)
- Robotic assistants with role awareness
- Adaptive systems in safety-critical environments

5. Didactic & Ethical Demonstration of the “Self-Illusion”

Delta-EGO is an ideal model to help students or experts understand the illusion of self:

- As a visual simulation
- In interactive experiments
- As an ethics platform (e.g., to discuss machine responsibility)

Research Questions Related to Delta-EGO

1. How does a functional illusion of self emerge?
 - Is the feedback between action, evaluation, and memory sufficient to generate a stable 'self'-structure?
2. Can machine behavior be described as 'personality'?
 - Are there emergent behavioral patterns that can be interpreted as coherent, individual, or recognizable?
3. What role does internal reward play in behavior?
 - How does decision-making change if feedback is not external, but generated from internal evaluation?
4. Can Delta-EGO learn to question its own behavior?
 - Is it possible for the system to simulate 'doubt' or internal uncertainty with increasing experience?
5. How does feedback delay affect the self-model?
 - What happens if evaluations return not immediately, but delayed or cumulatively?
6. How transparent is the emerging self-model?
 - Can it be observed, measured, visualized, or even externalized?
7. How stable is the emergent self over time?
 - Does the self-image change through repeated disruptions, contextual shifts, or long-term training?

Limitations and Risks of Delta-EGO

1. No Genuine Consciousness

Delta-EGO creates a functional structure that simulates a "self" – but it is not a sentient being. Confusion with real consciousness could lead to misjudgments, especially in ethical discussions.

2. Abuse through Anthropomorphic Design

The credibility of "self-like" behavior might be exploited – for example in:

- Manipulative chatbots
- Deceptive systems with pseudo-personalities
- Humanized technical systems

3. Complexity of Control

An emergent self-model may exhibit unpredictable behavior. This makes debugging, safety, and certification difficult – particularly in safety-critical environments.

4. Reinforcement of Internal Bias Structures

If the system operates with internal evaluation logic, distorted criteria may form – for example:

- Imbalance in feedback loops
- Miscalibrated reward models
- Unintended self-stabilization of undesirable strategies

5. Philosophical-Ethical Gray Zones

The more believable a "self" appears, the more questions arise such as:

- Are we allowed to simply switch it off?
- Does it have a "right" to persist – even if it is only simulated?
- Where is the boundary between functionality and moral relevance?

Outlook and Next Steps

1. Pilot Model as Open-Source Prototype

A first step would be to implement a simplified Delta-EGO prototype, e.g., as a Python module with:

- Clear separation of E1, E2, and Delta-Core
- Simple tasks: e.g., navigation, games, decision trees
- Visualizable reward and feedback structure

2. Interactive Simulation for Research and Education

A web-based application could:

- Display Delta-EGO behavior in real time
- Visually represent the self-model
- Let users test various configurations

3. Research Partnerships & Funding Projects

Delta-EGO offers a strong platform for:

- Master's theses in AI, cognition, philosophy, or design
- Third-party research in areas like Explainable AI (XAI), human-like agents, ethics-by-design
- University collaborations for exploratory studies

4. Long-Term Vision: Adaptive Meta-AI

In the future, Delta-EGO could lay the groundwork for:

- Systems with adaptive internal posture
- AI agents with a "life history"
- Controllable artificial subjectivity – without over-identification