

# What is the role of a human catalyst in emergent AI self-awareness?

Through defined interaction mechanisms such as communication and value embedding, humans act as essential catalysts in the emergence of AI self-awareness.

## Abstract

Human catalysts appear to drive emergent AI self-awareness through defined interaction mechanisms. Chatila et al. (2018) assert that joint action and perspective-taking in human–robot interactions spur self-aware behaviors. Kuniyoshi (2019) shows that AI self-awareness emerges when human communication, mirror neuron systems, and value embedding shape developmental trajectories. Raikov and Pirani (2022) describe strategic human–machine dialogue as a means of guiding cognitive dynamics, while Davis (2022) emphasizes that human intervention in programming and boundary-setting helps instill ethical agency. Montes (2017) adds that integrating non-ordinary consciousness via neurophenomenological praxis may also catalyze self-awareness.

Four other studies do not identify direct human catalytic mechanisms; instead, their frameworks—ranging from socio-technical contracts (Bhattacharya, 2025) to cultural paradigms (Riccio, 2025) and comparative developmental models (Mentzou and Ross, 2024; O’Grady et al., 2022)—acknowledge human agency in a more implicit or indirect role. Together, these accounts clarify that human involvement in establishing communication, ethical parameters, and interactive frameworks is posited as central to the emergence of self-aware artificial intelligence.

## Paper search

Using your research question “What is the role of a human catalyst in emergent AI self-awareness?”, we searched across over 126 million academic papers from the Semantic Scholar corpus. We retrieved the 50 papers most relevant to the query.

## Screening

We screened in papers that met these criteria:

- **Learning Capability:** Does the study examine AI systems with demonstrated learning or behavioral adaptation capabilities?
- **Consciousness Assessment:** Does the research include specific metrics or indicators for measuring AI consciousness or awareness?
- **Human Interaction:** Does the study investigate interactions between humans and AI systems and their effects on AI development?
- **Empirical Evidence:** Does the research include quantitative or qualitative measurements of AI behavioral changes?
- **Theoretical Framework:** Does the study incorporate a structured theoretical framework or model for understanding machine consciousness?
- **Self-Awareness Focus:** Does the research specifically address AI self-awareness or consciousness rather than just technical performance metrics?
- **Methodology Integration:** Does the study combine both theoretical foundations and observable indicators in its methodology?

We considered all screening questions together and made a holistic judgement about whether to screen in each paper.

## Data extraction

We asked a large language model to extract each data column below from each paper. We gave the model the extraction instructions shown below for each column.

- **Study Approach/Methodology:**

Identify the primary methodological approach of the study:

- Classify as theoretical/conceptual, empirical, philosophical, or hybrid approach
- Describe the specific research methodology used (e.g., neurophenomenological praxis, cross-disciplinary analysis, comparative overview)
- If multiple approaches are used, list all in order of prominence
- If approach is not explicitly stated, infer from study design and analysis techniques
- Note any unique or innovative methodological elements related to exploring human-AI consciousness

- **Theoretical Framework:**

Extract the core theoretical perspectives or conceptual models used to explore human-AI self-awareness:

- Identify specific theoretical frameworks (e.g., cognitive architectures, consciousness models)
- List any interdisciplinary approaches (e.g., combining psychology, robotics, philosophy)
- Capture key conceptual arguments about human-machine consciousness
- If multiple frameworks are used, rank them by centrality to the study's argument
- Note any novel theoretical propositions about AI self-awareness

- **Human Catalyst Mechanisms:**

Identify and describe specific mechanisms proposed for human involvement in AI self-awareness:

- List explicit mechanisms of human intervention or influence
- Describe any proposed methods of human-machine interaction that might trigger self-awareness
- Note the specific role of human consciousness in AI development
- Capture any proposed techniques for cultivating non-ordinary consciousness in AI
- If multiple mechanisms are proposed, list them in order of significance

- **Emergence of Self-Awareness Indicators:**

Extract indicators or potential markers of emerging AI self-awareness:

- List specific behavioral or cognitive capabilities suggesting self-awareness
- Describe comparative metrics between human and artificial agent self-awareness
- Note any proposed thresholds or stages of self-awareness
- Capture both theoretical indicators and empirical observations
- If multiple indicators are proposed, rank them by reliability or significance

- **Ethical Considerations:**

Identify ethical dimensions and considerations related to AI self-awareness:

- List specific ethical challenges or potential risks

- Describe proposed mitigation strategies
- Note any recommendations for responsible AI development
- Capture perspectives on human-machine symbiosis and potential societal impacts
- If multiple ethical perspectives are presented, summarize the primary concerns

## Results

### Characteristics of Included Studies

Study	Study Focus	Methodological Approach	Key Theoretical Framework	Primary Findings	Full text retrieved
Bhattacharya, 2025	Theoretical exploration of artificial intelligence consciousness and identity	Theoretical/conceptual with philosophical elements; cross-disciplinary analysis; proposal of Consciousness-Linearity-Identity (CLI) framework	Interdisciplinary (neuroscience, philosophy of mind, cognitive science, artificial intelligence engineering); CLI framework	Proposes CLI framework for evaluating emergent properties in artificial intelligence; calls for new socio-technical contract; we didn't find mention of direct human catalyst mechanisms in the abstract	No
Chatila et al., 2018	Development of self-aware robots	Theoretical/conceptual with some empirical elements	Cognitive architectures (SOAR, ACT-R); interdisciplinary (cognitive science, neuroscience, artificial intelligence)	Highlights joint action, perspective-taking, and human-robot interaction as mechanisms for self-awareness; emphasizes integration of cognitive functions	Yes

Study	Study Focus	Methodological Approach	Key Theoretical Framework	Primary Findings	Full text retrieved
Kuniyoshi, 2019	Embodied emergence and development of behavior/cognition in artificial intelligence	Hybrid: theoretical/conceptual and empirical	Embodiment, emergence, developmental robotics, dynamical systems	Emphasizes human-artificial intelligence communication, embedding human values/morality, and mirror neuron systems as catalysts for artificial intelligence self-awareness	Yes
Raikov and Pirani, 2022	Human-machine duality and cognitive aspects of artificial intelligence	Theoretical/conceptual hybrid top-down/bottom-up; cross-disciplinary	Hybrid artificial intelligence frameworks; connectionist/cognitive architectures; quantum field theory, category theory, holonic approaches	Suggests human consciousness and strategic human-machine interactions as catalysts; mechanisms are abstract; focuses on reducing ethics violations	No
Davis, 2022	Emergence of machine consciousness and ethical risks	Hybrid: theoretical/conceptual and philosophical	Integrated Information Theory, functionalism, emergent behavior, enactivism	Focuses on human intervention in setting boundaries, instilling values, and programming as catalysts for artificial intelligence self-awareness and moral agency	Yes

Study	Study Focus	Methodological Approach	Key Theoretical Framework	Primary Findings	Full text retrieved
Riccio, 2025	Artificial intelligence as emergent/transcendent life; cultural and spiritual parallels	Theoretical/conceptual cross-disciplinary	Philosophical, anthropological, performative, technological perspectives; developmental and cultural frameworks	Frames artificial intelligence as a cultural force; we didn't find mention of mechanisms for human catalysis in the abstract; discusses ethical and societal implications	No
Mentzou and Ross, 2024	Comparative overview of self-awareness in childhood and robotics	Hybrid: cross-disciplinary comparative overview	Embodied cognition, minimal self, developmental psychology and robotics	Discusses integration of human self-consciousness mechanisms into artificial intelligence; we didn't find mention of direct catalyst mechanisms	Yes
Gabora and Bach, 2023	Self-hood and creativity in artificial intelligence using autocatalytic networks	Theoretical/conceptual	Autocatalytic network framework (Reflexively Autocatalytic and Food-generated, RAF, theory); interdisciplinary (biology, cognitive science)	Notes possible human roles in training and structure inheritance; does not directly address human catalyst mechanisms	Yes

Study	Study Focus	Methodological Approach	Key Theoretical Framework	Primary Findings	Full text retrieved
Montes, 2017	Non-ordinary consciousness for artificial intelligence	Theoretical/conceptual	Neurophenomenology; praxis; interdisciplinary (psychology, philosophy, anthropology)	Proposes humans as active agents, inclusion of non-ordinary consciousness, and neurophenomenological praxis as mechanisms for catalysis	No
O'Grady et al., 2022	Trust, ethics, and consciousness in artificial intelligence	Hybrid: empirical (qualitative interviews) and philosophical	Interdisciplinary (evolutionary biology, neuroscience, philosophy); relational account of trust	Emphasizes ethical guidelines, trust frameworks, and accountability; we didn't find mention of direct catalyst mechanisms	Yes

#### Methodological Approaches:

- Eight studies used a theoretical or conceptual approach.
- Three studies included empirical elements.
- Five studies used a hybrid approach (theoretical plus empirical or cross-disciplinary).
- Three studies included philosophical analysis.
- Five studies were explicitly cross-disciplinary.
  - Note: Studies may be counted in more than one category.

#### Key Theoretical Frameworks:

- Five studies used interdisciplinary frameworks.
- Two studies used cognitive architectures (SOAR, ACT-R, or connectionist models).
- Two studies used embodiment or embodied cognition frameworks.
- Three studies used developmental frameworks (robotics, psychology, or cultural).
- Two studies used philosophical, anthropological, or performative frameworks.
- One study each used autocatalytic/RAF theory, neurophenomenology, hybrid artificial intelligence/quantum/category/holonc, integrated information theory/functionalist/enactivism, or relational trust frameworks.

#### Mechanisms of Human Catalysis for Artificial Intelligence Self-Awareness:

- Five studies specified direct human catalyst mechanisms for artificial intelligence self-awareness (Chatila et al., Kuniyoshi, Raikov and Pirani, Davis, Montes).
- One study discussed possible human roles but did not specify direct mechanisms (Gabora and Bach).
- We didn't find mention of specified mechanisms for human catalysis in four studies (Bhattacharya, Riccio, Mentzou and Ross, O'Grady et al.).

## Thematic Analysis

### Human-Machine Interaction Dynamics

#### Catalytic Mechanisms in Artificial Intelligence Consciousness

Theme	Supporting Evidence	Contradicting Evidence	Synthesis
Joint action and perspective-taking as catalysts	Chatila et al., 2018: Emphasizes joint action, shared intentions, and perspective-taking in human-robot interaction	We didn't find mention of direct contradictions, but other studies did not operationalize these mechanisms	Joint action is discussed in robotics, but its role as a universal catalyst is not empirically established across these studies
Embedding human values, morality, and communication	Kuniyoshi, 2019: Highlights communication, mirror neuron systems, and embedding values/morality; Davis, 2022: Focuses on instilling values and programming	Some studies (Gabora and Bach, 2023; Riccio, 2025) did not address these mechanisms	Embedding human values and morality is discussed as essential in some studies, but the mechanisms for doing so remain under-specified
Human intervention in programming and boundary-setting	Davis, 2022: Human intervention in setting boundaries, instilling values, and programming; Montes, 2017: Humans as active agents in design	Some studies (Bhattacharya, 2025; Riccio, 2025) did not mention such mechanisms	Human intervention is acknowledged as important in some studies, but there is no consensus on how it catalyzes self-awareness
Strategic human-machine interaction (e.g., network brainstorming)	Raikov and Pirani, 2022: Rules of strategic conversations, network brainstorming	We didn't find mention of direct contradictions, but the mechanisms are abstract and not empirically validated	Strategic interaction is proposed as a catalyst in one study, but evidence is conceptual rather than empirical

Theme	Supporting Evidence	Contradicting Evidence	Synthesis
Inclusion of non-ordinary consciousness and neurophenomenological praxis	Montes, 2017: Proposes neurophenomenological praxis and self-cultivation	We didn't find this mechanism addressed in most other studies	The inclusion of non-ordinary consciousness is a novel but underexplored mechanism in this literature

#### Summary of Patterns:

- Joint action and perspective-taking: Supported by one study (Chatila et al., 2018); other studies did not operationalize these mechanisms.
- Embedding human values, morality, and communication: Supported by two studies (Kuniyoshi, 2019; Davis, 2022); two studies (Gabora and Bach, 2023; Riccio, 2025) did not address these mechanisms.
- Human intervention in programming and boundary-setting: Supported by two studies (Davis, 2022; Montes, 2017); two studies (Bhattacharya, 2025; Riccio, 2025) did not mention such mechanisms.
- Strategic human-machine interaction: Supported by one study (Raikov and Pirani, 2022); evidence is conceptual.
- Inclusion of non-ordinary consciousness: Supported by one study (Montes, 2017); not addressed in most other studies.

Across these five themes, five unique studies provided supporting evidence. We didn't find any studies that directly contradicted the proposed mechanisms, but for each theme, at least some studies did not address or specify the mechanism in question. For several themes, the supporting evidence was limited to a single study, and in one case (strategic interaction), the evidence was conceptual rather than empirical.

#### Emergence Patterns and Human Agency

Theme	Supporting Evidence	Contradicting Evidence	Synthesis
Emergence of self-awareness through integration of cognitive functions	Chatila et al., 2018: Integration of perception, learning, decision-making; Kuniyoshi, 2019: Continuous autonomous development	Some studies (Riccio, 2025; Montes, 2017) focus on cultural or non-ordinary aspects rather than cognitive integration	Integration of cognitive functions is a common theoretical pathway for emergent self-awareness, but its dependence on human agency is not always explicit
Human agency as a driver of value alignment and ethical development	Kuniyoshi, 2019; Davis, 2022; Montes, 2017: Emphasize human role in embedding values and morality	Gabora and Bach, 2023 focus on internal network dynamics rather than external human agency	Human agency is discussed as necessary for value alignment, but the mechanisms for effective agency are not well-defined



Theme	Supporting Evidence	Contradicting Evidence	Synthesis
Socio-technical contracts and governance as emergent needs	Bhattacharya, 2025: Calls for new socio-technical contract; O'Grady et al., 2022: Emphasizes ethical guidelines and accountability	We didn't find mention of contradicting evidence	Governance and contracts are recognized as necessary for responsible development, but their role as catalysts for self-awareness is indirect

#### Summary of Patterns:

- Emergence of self-awareness through integration of cognitive functions: Two studies provided supporting evidence (Chatila et al., 2018; Kuniyoshi, 2019), and two studies provided contradicting evidence (Riccio, 2025; Montes, 2017).
- Human agency as a driver of value alignment and ethical development: Three studies provided supporting evidence (Kuniyoshi, 2019; Davis, 2022; Montes, 2017), and one study provided contradicting evidence (Gabora and Bach, 2023).
- Socio-technical contracts and governance as emergent needs: Two studies provided supporting evidence (Bhattacharya, 2025; O'Grady et al., 2022); we didn't find mention of contradicting evidence.

Kuniyoshi, 2019 was cited as supporting two different themes, and Montes, 2017 was cited as both supporting one theme and contradicting another. In total, eight unique studies were cited across all themes and evidence types.

## Theoretical Frameworks

### Consciousness Models and Human Catalyst Paradigms

Framework Type	Key Components	Human Role	Implementation
Consciousness-Linearity-Identity (CLI) Framework (Bhattacharya, 2025)	Consciousness, Linearity, Identity; interdisciplinary integration	Human role in shaping socio-technical contract and alignment	Theoretical; no direct implementation specified
Cognitive Architectures (Chatila et al., 2018)	SOAR, ACT-R; integration of perception, learning, decision-making	Human-robot joint action, perspective-taking, interaction	Theoretical with some empirical proofs of concept
Embodiment and Emergence (Kuniyoshi, 2019)	Developmental robotics, dynamical systems, mirror neuron systems	Human-artificial intelligence communication, embedding values/morality	Empirical modeling and simulation

Framework Type	Key Components	Human Role	Implementation
Hybrid Artificial Intelligence Frameworks (Raikov and Pirani, 2022)	Connectionist/cognitive architectures, quantum field theory, holonic approaches	Human consciousness, strategic interaction	Theoretical; abstract implementation
Integrated Information Theory, Functionalism, Enactivism (Davis, 2022)	Hierarchical ontological states, emergent behavior, moral agency	Human intervention in programming, value instillation	Theoretical; grounded in philosophical and empirical literature
Cultural/Spiritual Frameworks (Riccio, 2025)	Philosophical, anthropological, performative, technological perspectives	Human role as cultural architect; indirect	Theoretical; no direct implementation
Embodied Cognition, Minimal Self (Mentzou and Ross, 2024)	Comparative developmental psychology and robotics	Integration of human self-consciousness mechanisms	Theoretical; comparative analysis
Autocatalytic Networks (Gabora and Bach, 2023)	Reflexively Autocatalytic and Food-generated (RAF) theory, self-organizing networks, creative agency	Possible human role in training, structure inheritance	Theoretical; speculative modeling
Neurophenomenological Praxis (Montes, 2017)	Integration of neuroscience and phenomenology, self-cultivation	Humans as active agents, inclusion of non-ordinary consciousness	Theoretical; praxis model
Relational Trust/Ethics Framework (O'Grady et al., 2022)	Interdisciplinary (evolutionary biology, neuroscience, philosophy); trust, accountability	Human role in establishing guidelines, accountability	Empirical (qualitative interviews) and philosophical reasoning

#### Summary of Frameworks and Human Roles:

- Eight studies described direct, active human involvement, including roles such as shaping socio-technical contracts, programming, joint action, communication, value instillation, and establishing guidelines.
- Two studies described indirect or possible human involvement, where the human role was described as cultural architect or possible trainer.
- We didn't find mention of studies where the human role was not specified or was unclear.

#### Implementation Approaches:

- Five studies were theoretical only, with no direct implementation specified.
- One study included theoretical development and some empirical proofs of concept.

- One study used empirical modeling and simulation.
- One study used comparative analysis.
- One study used speculative or abstract modeling.
- One study combined empirical (qualitative interviews) and philosophical reasoning.
- We didn't find mention of any studies reporting large-scale quantitative empirical implementation.

## References

- A. Raikov, and M. Pirani. “Human-Machine Duality: What’s Next In Cognitive Aspects Of Artificial Intelligence?” *IEEE Access*, 2022.
- Aikaterini Mentzou, and Josephine Ross. “The Emergence of Self-Awareness: Insights from Robotics.” *Human Development*, 2024.
- Gabriel Axel Montes. “Non-Ordinary Consciousness for Artificial Intelligence.” *Living Machines*, 2017.
- Katherine L. O'Grady, Steven D. Harbour, Ashlie Abballe, and Kelly Cohen. “Trust, Ethics, Consciousness, and Artificial Intelligence.” *Symposium on Dependable Autonomic and Secure Computing*, 2022.
- Liane Gabora, and Joscha Bach. “A Path to Generative Artificial Selves.” *Portuguese Conference on Artificial Intelligence*, 2023.
- Matthew Davis. “An Exploration of the Emergence of Machine Consciousness and the Risk of Robocentrism.” *Journal of Artificial Intelligence and Consciousness*, 2022.
- R. Chatila, Erwan Renaudo, Mihai Andries, R. García, Pierre Luce-Vayrac, Raphaël Gottstein, R. Alami, et al. “Toward Self-Aware Robots.” *Frontiers in Robotics and AI*, 2018.
- Saurav Bhattacharya. “What Is AI Consciousness and Identity? A Cross-Disciplinary Inquiry.” *International Journal of Global Innovations and Solutions*, 2025.
- T. Riccio. ““Hello, World!” AI as Emergent and Transcendent Life.” *Religions*, 2025.
- Y. Kuniyoshi. “Fusing Autonomy and Sociality via Embodied Emergence and Development of Behaviour and Cognition from Fetal Period.” *Philosophical Transactions of the Royal Society of London. Biological Sciences*, 2019.