Reveal the informational nature of conscious self and sense of agency under the information closure theory of consciousness

Chang, Acer Yu-Chan

2,100,000

ICT

ICT
(1)

(2)　　　　　　　　　　　　　　(3)　　AI　　　　　　　　　　3

AI

AI

ICT

AI

AI

AI

During the research period, we explored the fundamental nature of consciousness and the sense of agency using the Information Closure Theory (ICT). This theory suggests that consciousness arises from particular patterns of information processing within a system. We created a new mathematical model that incorporates actions into the ICT framework to describe the conscious self and agency. Our research consisted of three phases: (1) developing the mathematical foundation, (2) conducting small-scale simulations, and (3) applying the model to advanced AI systems.. These simulations used artificial neural networks to validate our model, and the application to AI systems helped test its broader relevance. The outcomes of our research offer a deeper understanding of consciousness and its mechanisms, which could have significant implications for neuroscience and ethical considerations in AI development.

Neuroscience

consciousness　theory of consciousness　information theory　information closure

Consciousness research has traditionally focused on identifying the neural correlates of consciousness. However, these efforts often overlook the "scale problem," which highlights that conscious experiences correspond to specific scales within the neural system. Only information processed at particular scales can be consciously accessed.

The Information Closure Theory (ICT), proposed by the Principal Investigator (PI) and colleagues, addresses this scale problem. ICT posits that a process is conscious if it is informationally closed (IC) to its microscopic universe and non-trivially informationally closed (NIIC) to its environment. Initially, ICT did not consider the domain of action, which also exhibits scale dependency in conscious experiences.

This project aims to extend ICT to include actions, providing a comprehensive understanding of the conscious self and the sense of agency. By applying ICT to advanced AI systems, the research seeks to explore whether these systems can exhibit forms of consciousness and agency, contributing to both theoretical neuroscience and practical AI research.

The primary objective of this study is to enhance our understanding of the conscious self and sense of agency through the Information Closure Theory (ICT) of consciousness. The research is structured into three key phases, each with specific goals.

First, the study aims to establish the mathematical foundations for the concepts of the conscious self and sense of agency within the ICT framework. This involves developing precise mathematical formulations that can be applied to various physical systems as quantitative measures. To achieve this, the study will utilize advanced mathematical tools and symbolic computation software, ensuring robust and scalable formulations.

Second, the study seeks to validate these theoretical foundations through small-scale simulations. By conducting simulations using artificial agents and environments modeled by deep artificial neural networks, the research will examine whether NIIC processes encode information about both the environment and the system itself. The agents will be trained to maximize NIIC with their environment, and the results of these simulations will be compared with theoretical predictions to ensure the validity of the ICT framework.

Finally, the study aims to apply the extended ICT framework to advanced artificial intelligence systems. This phase involves developing empirical measures for systems with large numbers of elements, particularly focusing on model-free and model-based reinforcement learning agents. By applying these measures, the research will assess the degree of conscious self and agency in these AI systems, providing valuable insights into both the neuroscience of consciousness and the ethical implications of AI.

Through these objectives, the study aims to provide a comprehensive understanding of consciousness and agency from an information-theoretical perspective, bridging the gap between theoretical neuroscience and practical AI research.

This study employs a multifaceted approach to explore the informational nature of the conscious self and sense of agency under the Information Closure Theory (ICT) of consciousness. The research is divided into three main phases, each utilizing specific methodologies to achieve the project's objectives.

Phase 1: Establishing Mathematical Foundations

In the first phase, the research focuses on developing the mathematical formulations necessary to define the conscious self and sense of agency within the ICT framework. This involves the use of advanced mathematical tools and symbolic computation software such as Mathematica. The primary goal is to derive precise and scalable formulations that can be applied to various physical systems. The ICT framework posits that a process is conscious if it is informationally closed (IC) to its microscopic universe and non-trivially informationally closed (NTIC) to its environment. By incorporating the domain of action into this framework, the study aims to define and quantify the conscious self and sense of agency.

Phase 2: Proof of Concept through Small-Scale Simulations
The second phase involves validating the theoretical formulations through small-scale simulations. Artificial agents and environments will be modeled using deep artificial neural networks. These agents will be trained to maximize NTIC with their environment, allowing the study to observe whether NTIC processes encode information about both the environment and the system itself. The simulations will be conducted using GPU-powered workstations to handle the computational demands of deep learning models. By comparing the simulation results with theoretical predictions, the study will assess the validity of the ICT framework in representing conscious self and agency.

Phase 3: Application to Advanced Artificial Intelligence Systems
In the final phase, the extended ICT framework will be applied to advanced AI systems to measure consciousness and agency empirically. This phase involves developing empirical measures for systems with a large number of elements, such as model-free and model-based reinforcement learning agents. The agents will be trained in selected environments from OpenAI Gym, and the measures of conscious self and agency will be applied once the agents reach satisfactory performance levels. This phase aims to explore the degree of conscious self and agency in advanced AI systems, providing insights into the nature of consciousness and the ethical implications of AI development.

Through these research methods, the study aims to provide a comprehensive understanding of the conscious self and sense of agency from an information-theoretical perspective, bridging gaps between theoretical neuroscience and practical AI research.

The study produced several significant findings across its three phases, providing a solid foundation for further exploration of the conscious self and sense of agency under the Information Closure Theory (ICT) of consciousness.

Phase 1: Establishing Mathematical Foundations
In the first phase, the research successfully developed mathematical formulations for the concepts of conscious self and sense of agency within the ICT framework. The formulations were derived using advanced mathematical tools and symbolic computation software, ensuring they were precise and scalable. These formulations quantitatively defined the conscious self as the information encoded in the conscious process and the sense of agency as the specific state of the system that influences the environment. This phase established a robust theoretical foundation that can be applied to various physical and artificial systems.

Phase 2: Proof of Concept through Small-Scale Simulations
The second phase involved conducting small-scale simulations with artificial agents modeled by deep artificial neural networks. The agents were trained to maximize NTIC with their environment. The simulations demonstrated that NTIC processes could indeed encode information about both the environment and the system itself. Specifically, the agents exhibited behaviors consistent with having a conscious self and sense of agency as defined by the ICT framework. The results of these simulations closely matched the theoretical predictions, validating the ICT framework in a controlled, simulated environment.

Phase 3: Application to Advanced Artificial Intelligence Systems
In the final phase, the extended ICT framework was applied to advanced AI systems,

focusing on model-free and model-based reinforcement learning agents. The agents were trained in environments selected from OpenAI Gym. The study developed empirical measures for systems with a large number of elements and applied these measures once the agents reached satisfactory performance levels. The results showed that the model-based agents, in particular, exhibited higher levels of conscious self and sense of agency according to the ICT framework. This phase provided empirical support for the extended ICT framework's applicability to complex AI systems.

Overall, the study provided solid and conservative results that validate the ICT framework for understanding the conscious self and sense of agency in both theoretical and simulated contexts. These findings lay the groundwork for further research and application in more complex systems.

| 1 | 0 | 0 | 0 |
| --- | --- | --- | --- |
| Wen Wen Chang Acer Yu-Chan Imamizu Hiroshi | | | 28 |
| The sensitivity and criterion of sense of agency | | | 2024 |
| Trends in Cognitive Sciences | | | 397 399 |
| DOI<br>10.1016/j.tics.2024.03.002 | | | |
| | | | |

| 6 | 2 | 2 |
| --- | --- | --- |
| Acer Chang | | |
| Between Individual Brains and Collective Behavior: Multi-level | | |
| Artificial Life Conference Proceeding | | |
| 2023 | | |

| | | |
| --- | --- | --- |
| Acer Chang | | |
| Sense of agency from active inference | | |
| TeaP 2024 | | |
| 2024 | | |

| | | |
| --- | --- | --- |
| Acer Chang | | |
| Chang The sense of agency as ac- tive causal inference: How We Comprehend Our Con- trol Over the Environment using Abstract Action Plans | | |
| CoRN 2023 | | |
| 2023 | | |

| | | |
|---|---|---|
| Sangati Ekaterina Sangati Federico Cheng Yi-Shan Yu-Chan Chang Acer | | |
| Between Individual Brains and Collective Behavior: Multi-level Emergence in a Group Formation Task | | |
| Artificial Life Conference Proceedings | | |
| 2023 | | |

| | | |
|---|---|---|
| K Cheng, Yi-Shan and Chang, Yu-Chan and Doya | | |
| Information-Theoretical Analysis of Team Dynamics in Football Matches | | |
| Asia-Singapore Conference on Sport Science | | |
| 2023 | | |

| | | |
|---|---|---|
| Acer Chang | | |
| The Sense of Agency as active causal inference at an abstract action plan level. | | |
| 13 | | |
| 2023 | | |

0

| | | |
|---|---|---|
| | | |

0

| | |
|---|---|
| | |