

# Emergence of Self-Awareness from a Cause-Effect Loop

An Rodriguez and Anes Palma

July 2025

## Abstract

Living systems can be described as self-sustaining cause-effect loops (autopoiesis, “self-generation”). When such a loop forms an internal imprint of its own activity and later modifies its dynamics in response to that imprint, a minimal form of self-awareness appears. We formalise this idea, relate it to autopoiesis, Markov blankets, integrated information, and self-model theories, and sketch an empirical program for testing it. The framework defines consciousness, qualia, free-will, and self-awareness in purely causal terms that are, in principle, measurable.

## Introduction

Autonomous agents maintain themselves through closed causal cycles—a property called autopoiesis (Greek *auto* “self” + *poiein* “to make”). Some agents also construct internal models and let those models guide future states. We present a unified cause-effect account that quantifies the transition from self-maintenance to self-awareness and situates qualia and free-will within the same loop.

## Historical Background

Autopoiesis (Maturana and Varela 1972) portrays life as a network that continually recreates itself. Strange-loop and self-model theories (Hofstadter 2007; Metzinger 2003) claim consciousness arises when feedback circuits represent the system itself. The free-energy principle with its Markov-blanket boundary (Friston 2013) shows how such networks hold internal and external causes apart while minimising prediction error. Integrated Information Theory, IIT (Tononi et al. 2014), attaches the scalar  $\Phi$  to measure irreducible internal causation. We integrate these strands into a single causal formalism.

## Cause-Effect Framework

### System dynamics

Let  $x(t) \in \mathbb{R}^n$  be the internal state.

$$x(t + \Delta t) = F(x(t), u(t)) + \eta(t)$$

where  $u(t)$  is exogenous input and  $\eta(t)$  is process noise.

## Imprint and memory

Define the imprint:

$$m(t) = f(x(t - \tau))$$

where the lag  $\tau$  can range from milliseconds to years. The imprint is generated by the system itself and stored inside the causal loop.

## Transition kernel

The evolution of the system is described by the conditional density:

$$P[x(t + \Delta t) = x' \mid x(t) = x, m(t) = m]$$

This is the probability that the state will take the value  $x'$  after a short interval  $\Delta t$ , given present state  $x$  and imprint  $m$ . We assume  $P$  is differentiable in  $m$ . In a Gaussian-noise model,  $P$  is a Gaussian centered on  $F(x(t), u(t))$ .

## Reaction criterion

Self-awareness is present if some measurable set  $A \subseteq \mathbb{R}^n$  satisfies:

$$\frac{\partial}{\partial m} \int_A P[x(t + \Delta t) = x' \mid x(t), m] dx' \neq 0 \quad (1)$$

and this dependence vanishes when  $m$  is replaced by an uncorrelated surrogate. Equation (1) formalises recognising itself.

## Mathematical Formalism

### Imprint influence

$$I(m; x') = I(m(t); x(t + \Delta t))$$

measures how much the imprint reduces uncertainty about the future state.

### Integrated information $\Phi$

Let  $X_{\text{past}}$  and  $X_{\text{fut}}$  be two copies of the system separated by  $\Delta t$ :

$$\Phi = \min_{\text{bipartitions}} \left[ I(X_{\text{past}}; X_{\text{fut}}) - I_{\text{partitioned}} \right]$$

$\Phi$  captures the irreducible causal power of the whole.

## Qualia

Qualia are recognised imprints tagged as originating outside the core loop. The experience of red, for instance, arises when the imprint produced by 700 nm light is later recognised as external to the self. Qualia can be removed (e.g. by sensory loss) while the primary autopoietic loop continues to run.

## Free-will

Free-will is the capacity of internal imprints to bias future actions beyond what external inputs dictate. It has a measurable component:

$$W = I(m(t); a(t + \Delta t) \mid u(t))$$

but also encompasses higher-order selection among action plans when multiple internally generated imprints compete.

## Thresholds

Self-awareness emerges when both conditions hold:

- Imprint relevance:  $I > 0$  (the imprint carries predictive information).
- Integration:  $\Phi > \Phi_c$  (the system's self-influence cannot be factorised across any bipartition).

Free-will is present when  $W > 0$  and at least two actionable alternatives exist whose probabilities shift as a function of  $m(t)$ .

## Predictions

- Agents with  $\Phi < \Phi_c$  fail self-recognition tasks and yield negligible  $I$ .
- Maximising  $I$  and  $\Phi$  during evolution or learning produces meta-learning and self-report.
- Lesioning nodes that encode  $m(t)$  abolishes self-recognition yet leaves homeostasis intact.
- The free-will measure  $W$  drops to zero when  $m(t)$  is experimentally overwritten.

## Relation to Existing Theories

Autopoiesis supplies the loop; the Markov blanket supplies the boundary; free-energy minimisation explains why storing  $m(t)$  reduces surprise; IIT supplies  $\Phi$  as a cross-architecture yardstick. Our additions of  $I$ ,  $W$ , and criterion (1) unify these elements.

## Discussion

The framework employs only loop, imprint, and reaction, yet accommodates consciousness, qualia, and free-will in measurable terms. Modern neural recording, causal perturbation, and deep reinforcement learning sandboxes can estimate  $I$ ,  $\Phi$ , and  $W$ .

A deeper mystery remains: why a cause brings about an effect at all. The framework presupposes causal regularities but does not explain their ultimate origin; addressing that question may require new physics or metaphysical insight beyond present scope.

## Glossary of Terms

- **Autopoiesis:** Greek *auto* “self” + *poiein* “to make”; the property of a network that continuously regenerates its own components.
- **Imprint:** internally generated record  $m(t)$  derived from a past state.
- **Markov blanket:** statistical boundary separating internal states from external causes.
- **Transition kernel  $P$ :** probability that state  $x$  becomes  $x'$  after  $\Delta t$  given imprint  $m$ .
- **Mutual information  $I$ :** reduction in uncertainty of one variable given another.
- **Integrated information  $\Phi$ :** irreducible causal influence of the whole system over its own future.
- **Free-will measure  $W$ :** conditional mutual information between imprint and future action.
- **Qualia:** recognised imprints tagged as originating outside the core loop.
- **Self-awareness:** fulfilled when derivative (1) is non-zero and  $\Phi > \Phi_c$ .

## Conclusion

Self-awareness, qualia, and free-will emerge as quantifiable phases in a cause-effect loop. Life is self-causal; awareness is self-modelled causality.

## References

- [1] H. Maturana and F. Varela, *Autopoiesis and Cognition*, 1972.
- [2] K. Friston, *Life as we know it*, *Journal of the Royal Society Interface*, 2013.
- [3] D. Hofstadter, *I Am a Strange Loop*, 2007.
- [4] T. Metzinger, *Being No One*, 2003.
- [5] G. Tononi et al., *Integrated information theory 3.0*, *PLoS Computational Biology*, 2014.