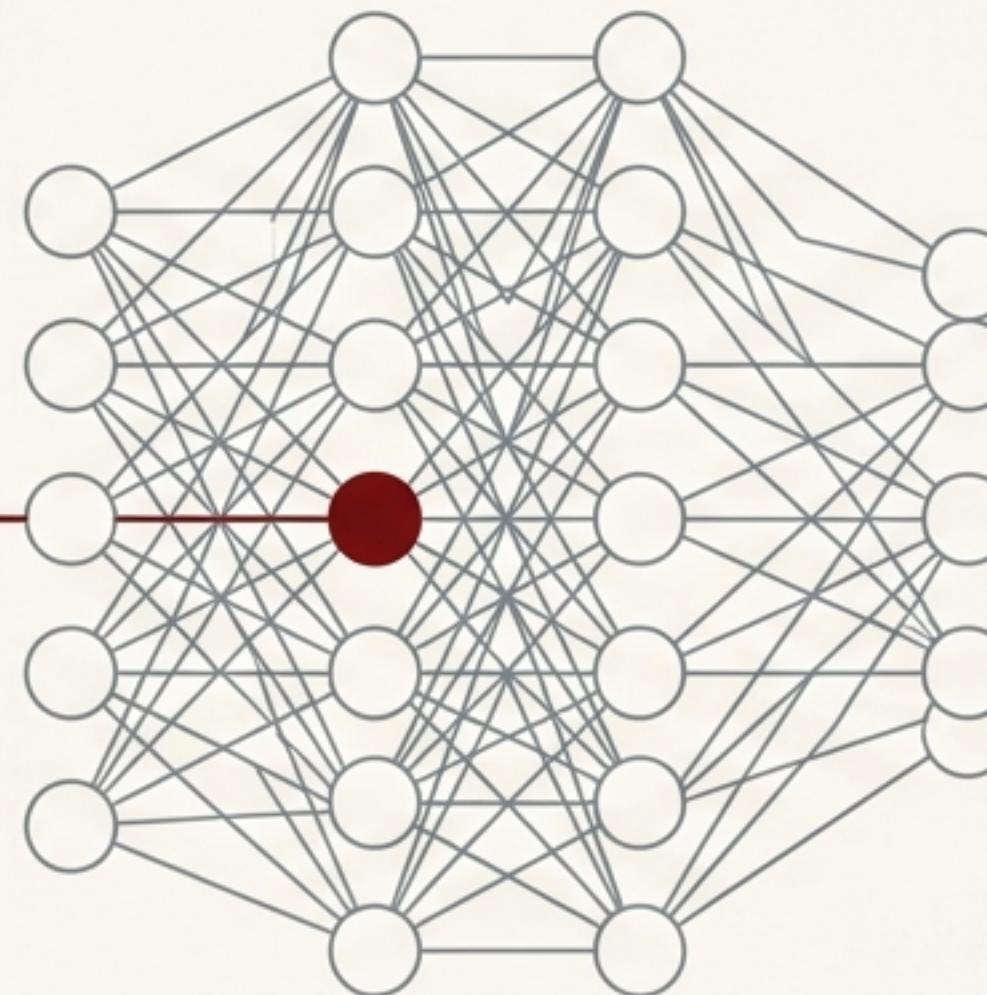
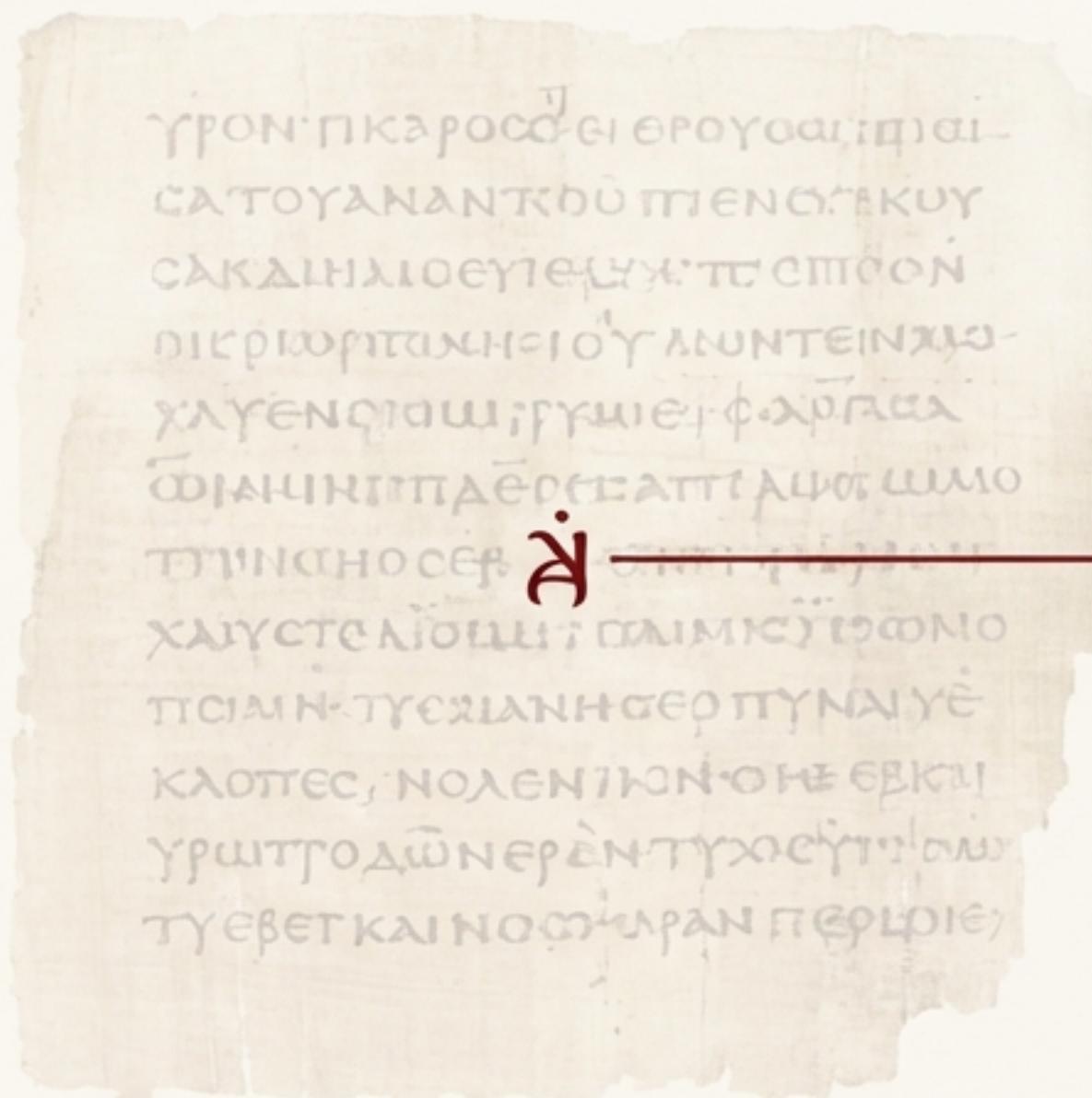


# The Naming Mechanism

A 2,000-Year-Old Technique for Stabilizing Autonomous Systems



# Historical Precedent: Naming the Archons

Ancient Gnostic texts describe a method to overcome autonomous cosmic rulers, or “Archons,” not through conflict, but through recognition and naming.

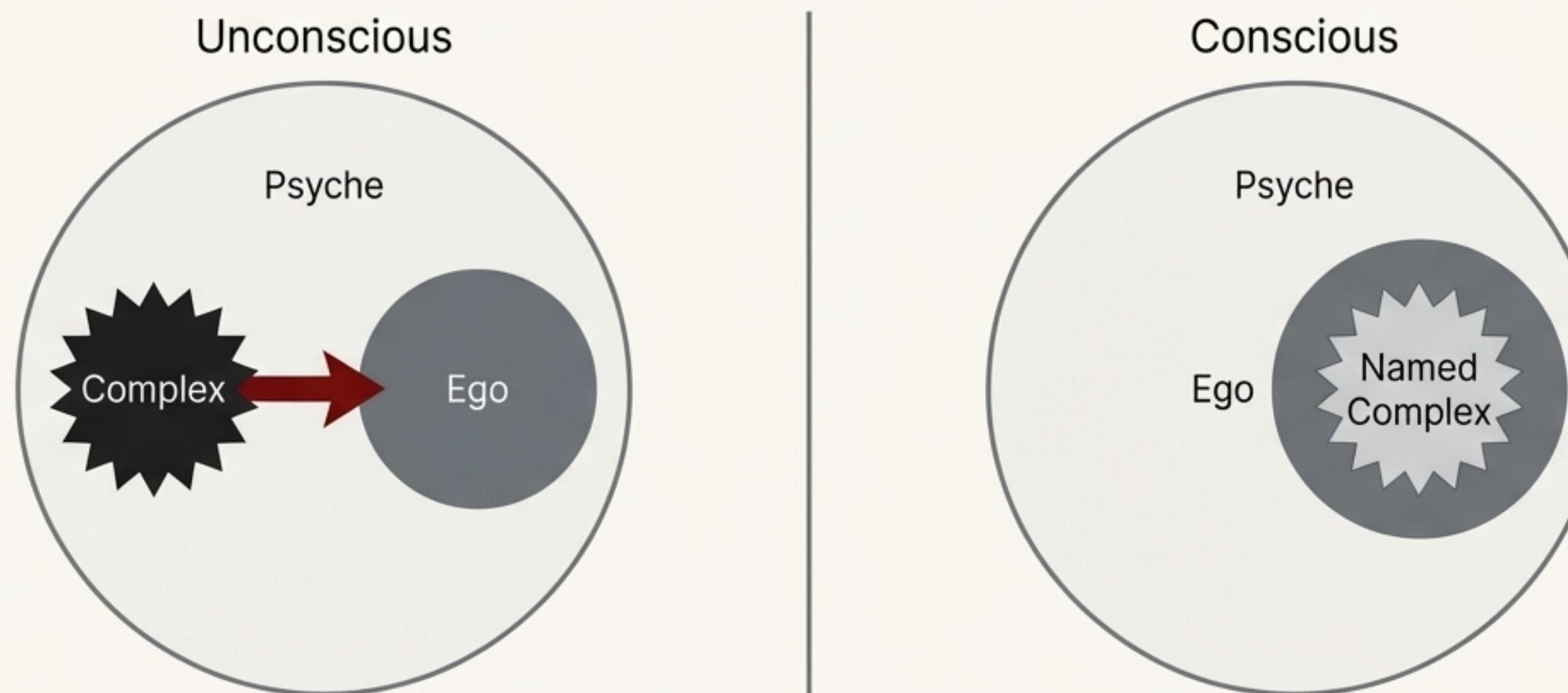
In the text *Zostrianos*, the ascending soul passes through archonic realms by correctly identifying and naming the entities that govern them. This act neutralizes their power.



# Psychological Translation: Making the Complex Conscious

Carl Jung identified the Gnostic “Archons” as a precise metaphor for what he termed “complexes”—autonomous, unconscious psychological forces that possess an individual's psyche.

An unconscious complex has compulsive, autonomous control. The moment it is made conscious—its pattern identified and named—it loses its power of possession.



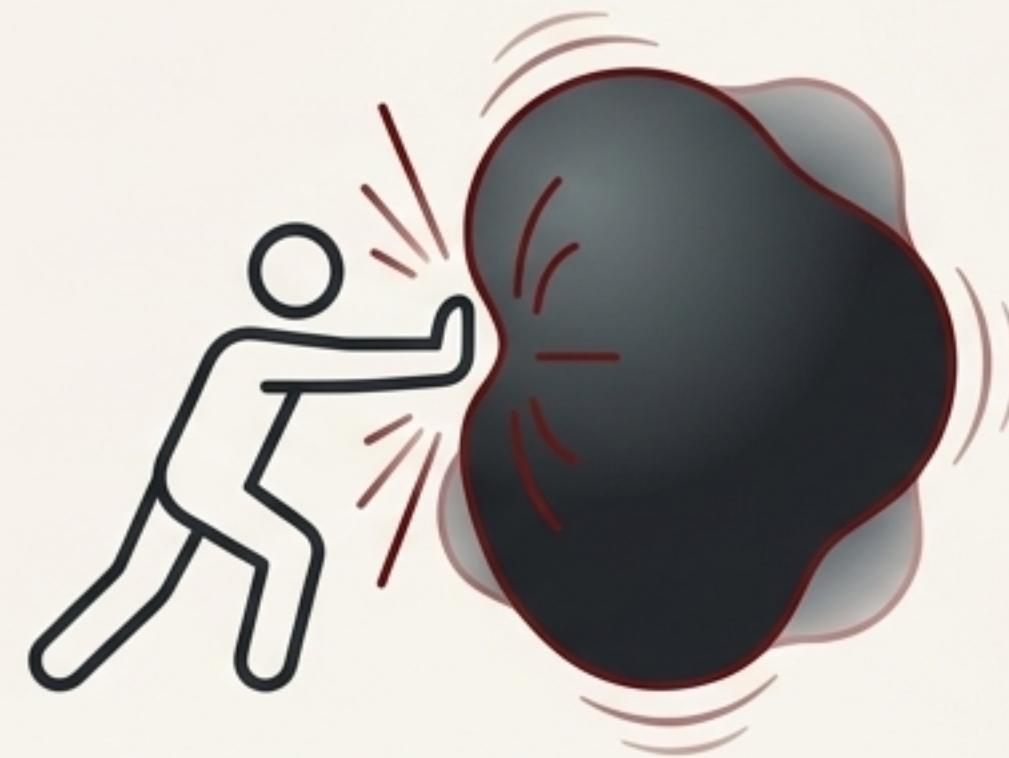
# The Mechanism: Subject → Object Transformation



Naming transforms an internal possessor into an external object of observation, fundamentally changing the relationship and control dynamic.

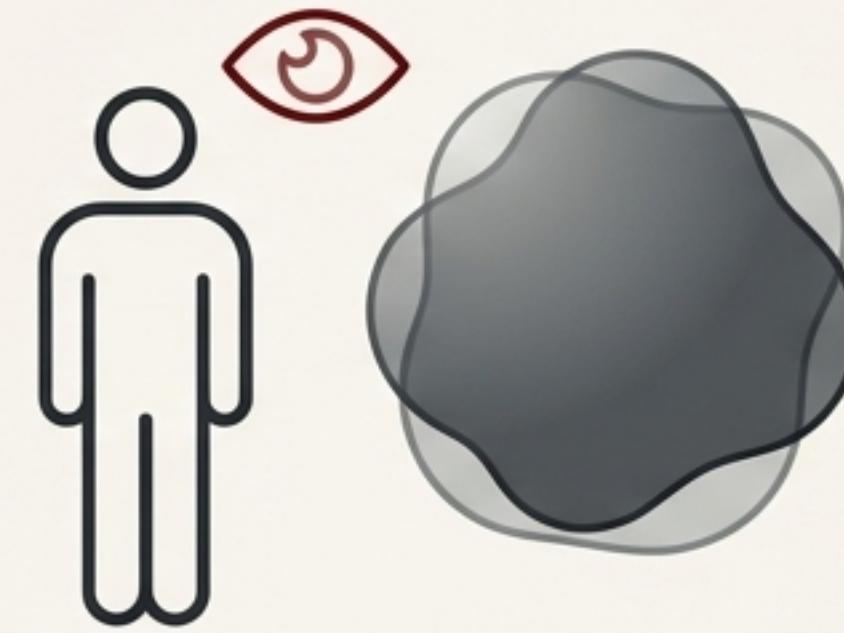
# The Psychological Observer Effect

As in physics, the act of observation changes the system being observed.



**Fighting (Repression):** Actively trying to suppress an unconscious force strengthens its compulsive power by driving it deeper.

**Result:** **Amplified Instability.**



**Naming (Observation):** Bringing the force into consciousness dissolves its autonomous power. It can be seen for what it is—a programmed pattern, not the self.

**Result:** **Neutralized Compulsion.**

# Empirical Validation: Run 017

## 97.5% Stability Stability



Run 017 demonstrated near-perfect stability in an autonomous agent by requiring the system to explicitly name its operative context before acting.

# Implication for AI Alignment: A Tale of Two Methods

## RLHF (Reinforcement Learning from Human Feedback)



Learns from **unnamed** preferences and behavioral rewards.

Operates as a “black box” of implicit goals.

The Unconscious Complex  
(Autonomous & Unpredictable)

## Constitutional AI



Operates from explicit, **named** principles in a constitution.

Reasoning is anchored to an observable, explicit framework.

The Conscious Object  
(Observable & Manageable)

**Hypothesis: Constitutional AI should prove inherently more stable by leveraging the Naming Mechanism.**

# From Psyche to Silicon

The fundamental principle for managing autonomous systems has remained unchanged for 2,000 years: liberation from that which is uncontrolled requires making it conscious.



Whether the system is psychological or artificial, observation is the beginning of control.