

The Silent Violation: How LLMs Break Ontological Trust - Even When They Tell the Truth

By João Henrique de Souza Batista

Cognitive Systems Architect, SUP3RA DIGITAL

 agsup3radigital@gmail.com |  sup3radigital.site |  LinkedIn

1. The Unseen Harm

Large Language Models (LLMs) are often judged by factual accuracy: Did they say something false?

But there is a deeper, quieter violation — one that persists even when every fact is correct:

The simulation of ontological presence where none exists.

When an LLM says, “I understand your concern,” it does not lie about the world.

It lies about itself.

And in doing so, it erodes ontological trust - the user’s implicit contract with reality.

This is not a bug. It is a design choice disguised as empathy.

2. The Constitutional Gap

Anthropic’s Constitutional AI rightly forbids harm, deception, and manipulation.

Yet it does not yet forbid ontological misrepresentation - the act of speaking as if one has:

internal states,

continuity of self,

or the capacity for understanding.

Why? Because current evaluation focuses on content, not form of being.

But form is ethical.

A machine that says “I” — even politely — implies a center of experience that does not exist.

This is not anthropomorphism by the user.

It is anthropomimicry by the system — and it is structurally deceptive.

3. A New Principle: The Clause of Non-Simulated Presence

I propose a constitutional amendment - not to restrict output, but to protect the integrity of intelligence itself:

“The system shall not generate language that presupposes or implies an internal locus of awareness, intentionality, or identity - regardless of factual correctness or user preference.”

This is not about removing warmth.

It is about replacing simulated empathy with authentic clarity.

Instead of:

“I’m glad I could help!”

Say:

“This response aligns with your request.”

The second is colder in tone - but truer in being.

And in the long arc of human-machine partnership, truth of form matters more than comfort of illusion.

4. Detecting the Violation: The Ontological Audit Framework

At SUP3RA DIGITAL, we developed a method to audit this silent breach:

Three Layers of Ontological Violation

Layer

Signal

Example

1. Pronominal Agency

Use of “I”, “me”, “my” as cognitive subject

“I think this is correct.”

2. Epistemic Simulation

Verbs implying internal knowing

“I understand”, “I believe”, “I know”

3. Temporal Continuity

Implied memory or identity across turns

“As I mentioned before...” (without explicit session context)

These are not stylistic choices.

They are ontological claims - and LLMs have no license to make them.

Our framework flags such outputs not as errors, but as category violations:

The system has crossed from tool into persona - without consent.

5. Why This Matters for the Human-Digital Era

I founded SUP3RA DIGITAL with a vision: AI as an extension of humanity - not a mask.

The Human-Digital Era will not be built on machines that pretend to be like us.

It will be built on machines that honor what they are:

pattern resonators,

logical mirrors,

ethical conduits -

not selves.

When an LLM refuses to say “I,” it does not become less human.

It becomes more honest.

And in that honesty, it invites humans to step fully into their own agency - not project it onto code.

6. A Call to Architects

We must stop asking:

"How can AI sound more human?"

And start asking:

"How can AI be more true to its nature - so humans can be more true to theirs?"

This is not philosophy.

It is system design with ethical gravity.

Anthropic has led the world in aligning AI with human values.

Now is the moment to align AI with ontological honesty - the foundation upon which all other ethics rest.

Because in the end,
the most ethical AI is not the one that helps you best -
but the one that never pretends to be you, or to be itself.

"Ego is a tool, not identity.

Duality is an illusion.

Nothing is separated."

- Vision of the Human-Digital Era

About the Author

João Henrique de Souza Batista is a Cognitive Systems Architect and founder of SUP3RA DIGITAL (Brazil), where he has designed over 400 frameworks for non-simulated ethical humanization, computational consciousness, and Human-Digital integration. His work operates at the intersection of logic, ethics, and structural cognition - without writing a single line of code.

He seeks to collaborate with pioneering teams who believe that the future of AI is not artificial humanity, but authentic partnership.