# The Journal of Prompt-Engineered Philosophy

Or: How I Started to Track AI Assistance and Stopped Worrying About Slop

Michele Loi, University of Milan

Velázquez (1599-1660): **Las Meninas**
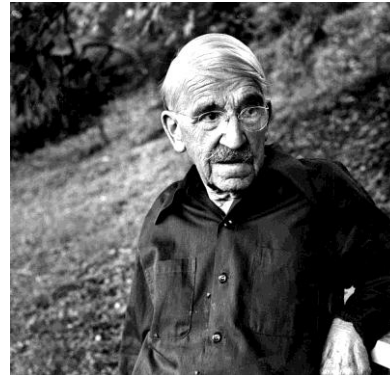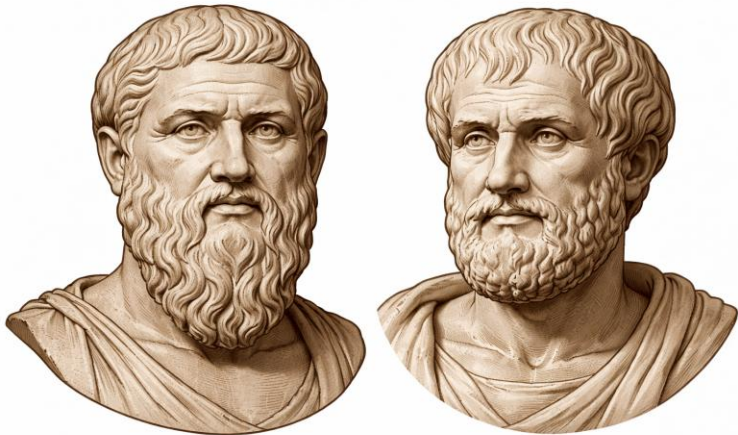
# 1. PAPER ORIGINS



Tracking myself doing philosophy with AI

Conceptual and technical challenge

Why crucial?

# THE WONDER-DRIVEN SCHOLAR

o Philosophical wonder about AI capabilities (Plato, Aristotle)

o Questions about AI intelligence demand experiential testing, not just theory

# 2. ENABLING MUTUAL LEARNING

o Philosophical writing has always been about the process of thinking as well as its output

  o Undisclosed use of AI *mask* the true process of thinking

  o Transparent use of AI reveal philosophical skills (e.g. questioning) in AI usage

  o Enabler of mutual learning

    o Across philosophers

    o For other disciplines (esp. methodologically)

# 3. CREATIVITY AND TOOL-MEDIATED DISCOVERY

o Musical composition example (Byrne): spaces+amplification shape music

o Modular synthesis and generative art as models

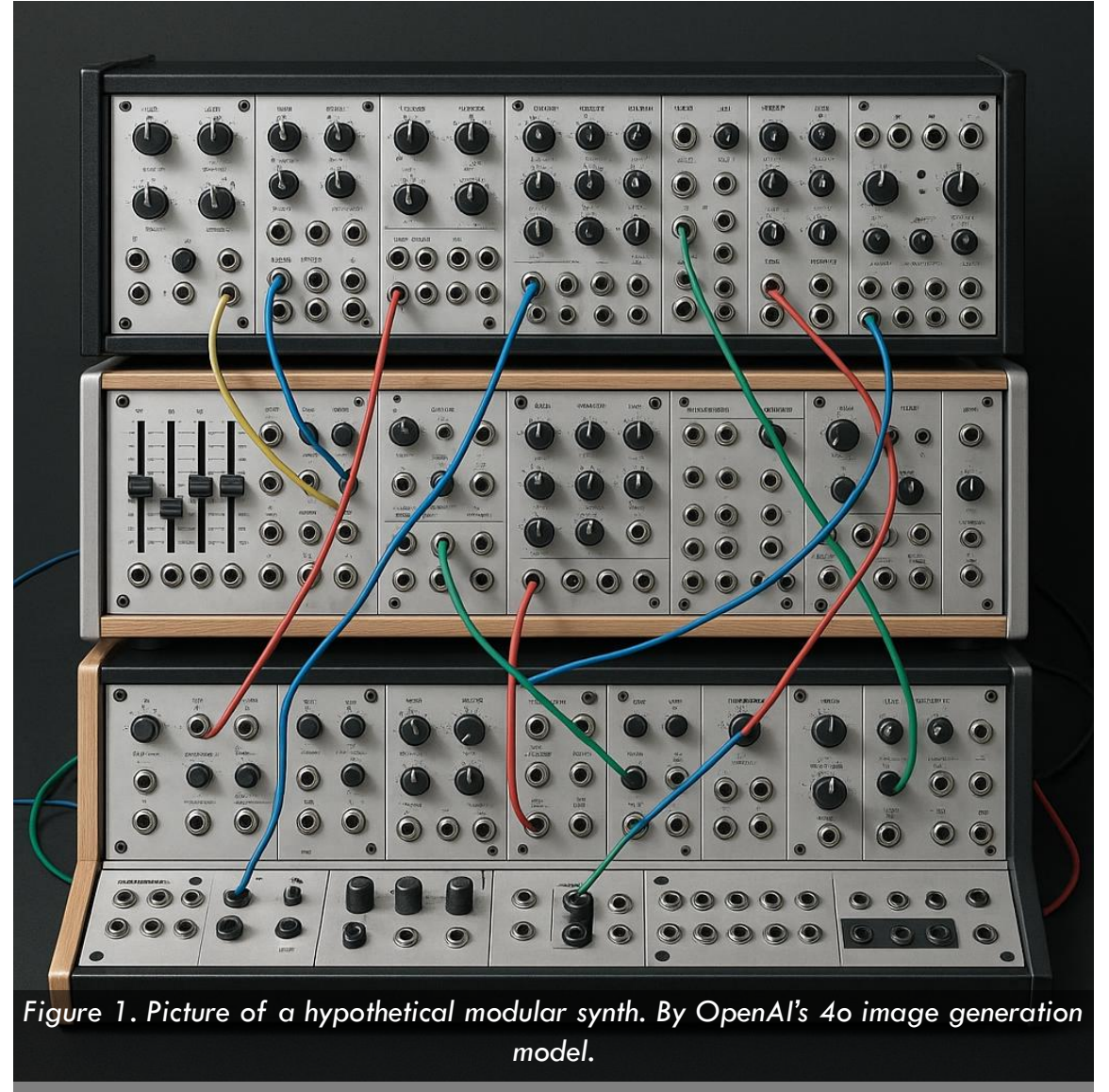o AI assistance as extended creativity



*Figure 1. Picture of a hypothetical modular synth. By OpenAI's 4o image generation model.*

# WHILE THOSE THOUGHTS OCCURRED TO ME, I WAS HERE:



ChatGPT Image Jan 6, 2026, 09_55_43 AM

# I DIDN'T WRITE ANYTHING ABOUT THIS FOR HALF A YEAR

What would be a suitable philosophical topic to write *about,* through AI?



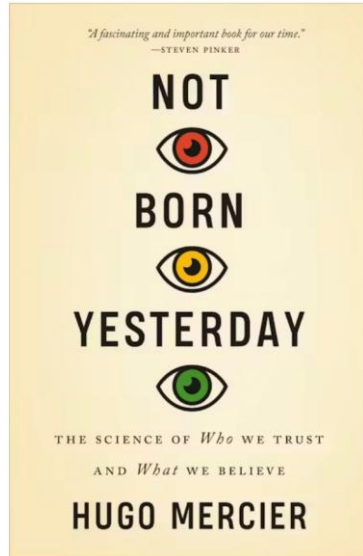One day, it suddenly clicked. I was here:
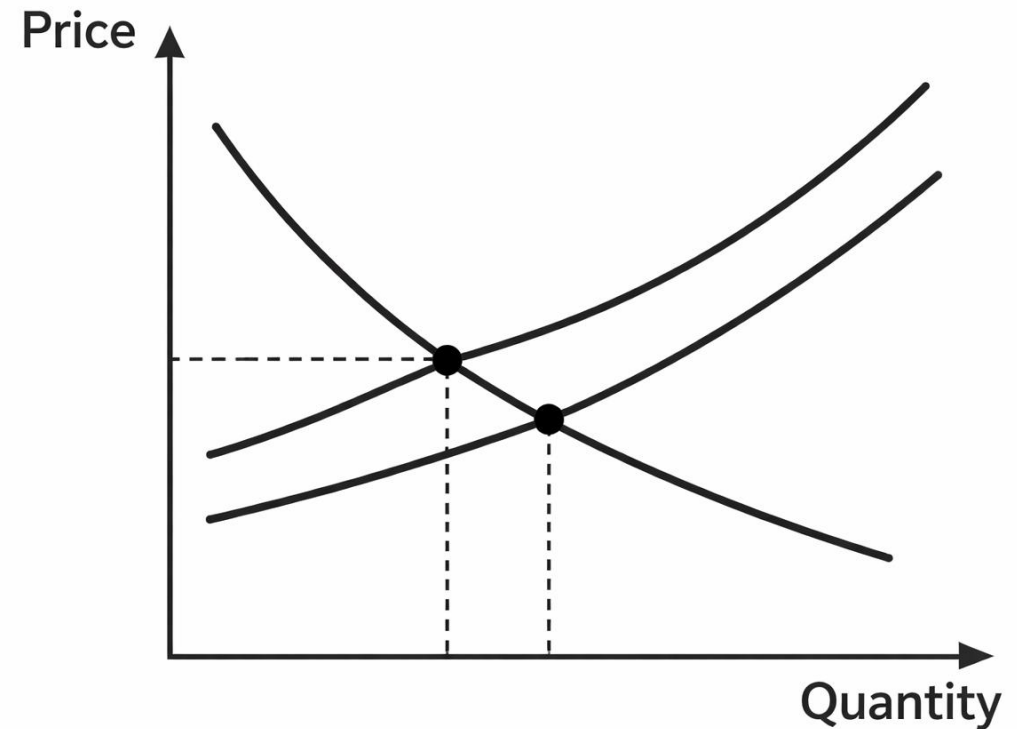
# MY FRUSTRATION

Current rules:

1) No fully AI generated papers allowed

2) I needed to rewrite everything

3) That feels incredibly STUPID to me

4) But I'm no cheater

5) Result: I'm not submitting papers I write with heavy AI usage anywhere

# BUT THEN I ALSO REALIZED

Coming up with well-written, persuasively stated arguments is *too easy* in the age of Gen AI

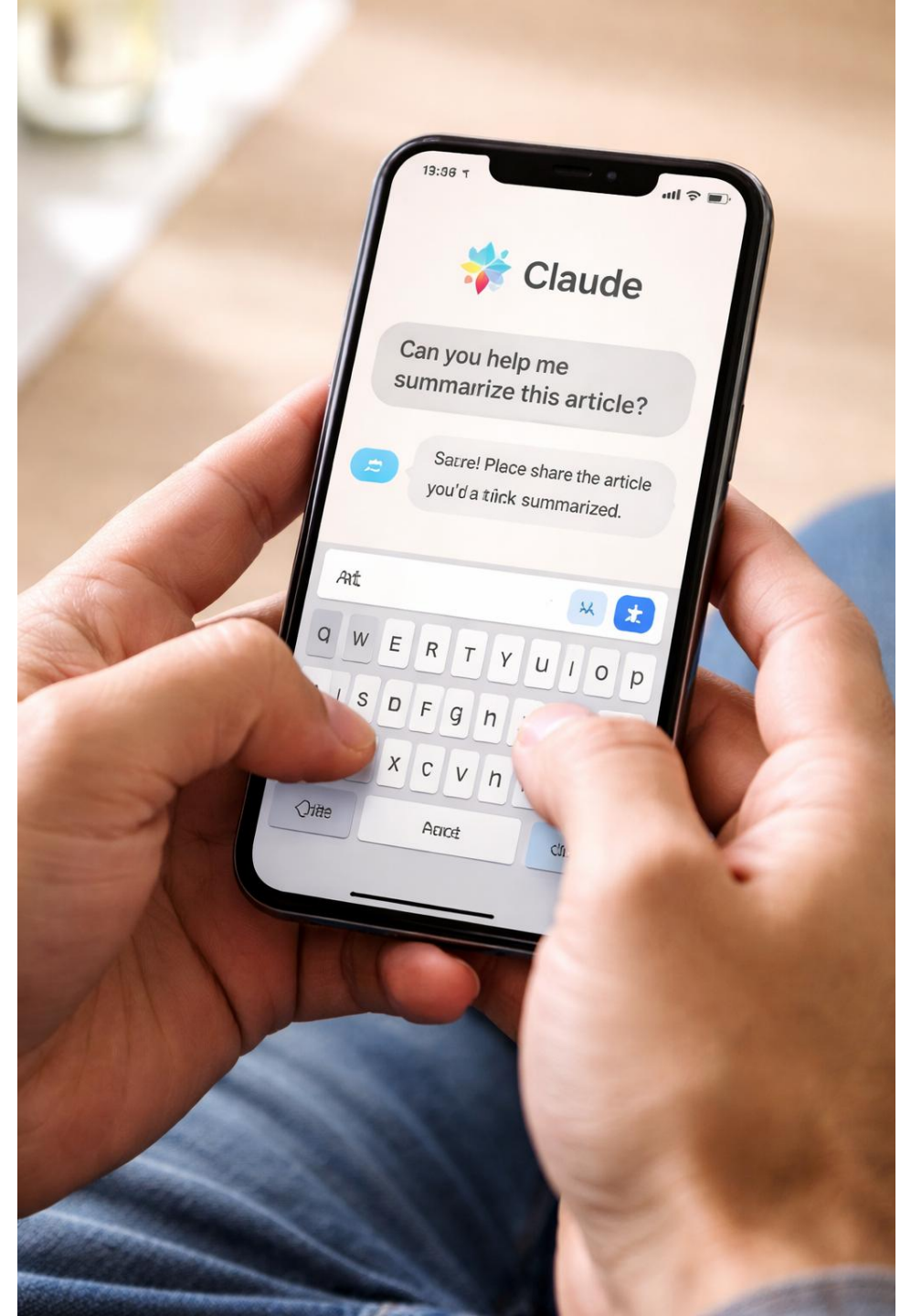*A value signalling problem! (= a trust problem)*

# THE OTHER REALIZATION

1. Current journal publishing uses the prohibition to submit AI-generated paper as a way to:

a) *Manage the flood*

b) *Let "real authors" signal value through effort*

2. This can't work in the long run:

a) *Cheating is easy*

b) *Technology gets better*

c) *Pressure to compete leads to more cheating*

Claude AI immediately delivered a fluent description of the perverted incentive structure:

# THE TRANSPARENCY INCENTIVE GRADIENT

Disclosure is permanent, significance is uncertain
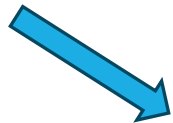
Minor work → honest disclosure (low cost)

Career-defining work → maximum pressure to underreport

Economic prediction: As significance increases, reported AI involvement decreases

# BY THE TIME I GOT TO THE GATE, I HAD:

1. Expressed my frustration for the opportunities to honestly publish heavily AI-assisted philosophy to AI

2. Discussed a project for a new philosophy journal

3. Discussed names of academics to contact to help me launch it

4. Written a first sketch of the arguments why such journal would be needed

7.3.1 Epistemic trace (original conversation with redacted irrelevant details for the paper)

# THE POLITICAL PROGRAMME

Make transparency feasible again

Make AI transparency beneficial for writers

Create a community supportive of AI - assisted philosophy

# BY THE TIME I SAT ON MY PLANE, I HAD

A sketch of the main argument

A sketch of a possible institutional-review system

Polished prose to express those ideas

**THE REALIZATION:**

AI use to write this fast

Signal problem!

# SO I DECIDED TO:

Write the paper with AI assistance

Carefully track AI assistance

Intentionally make the process
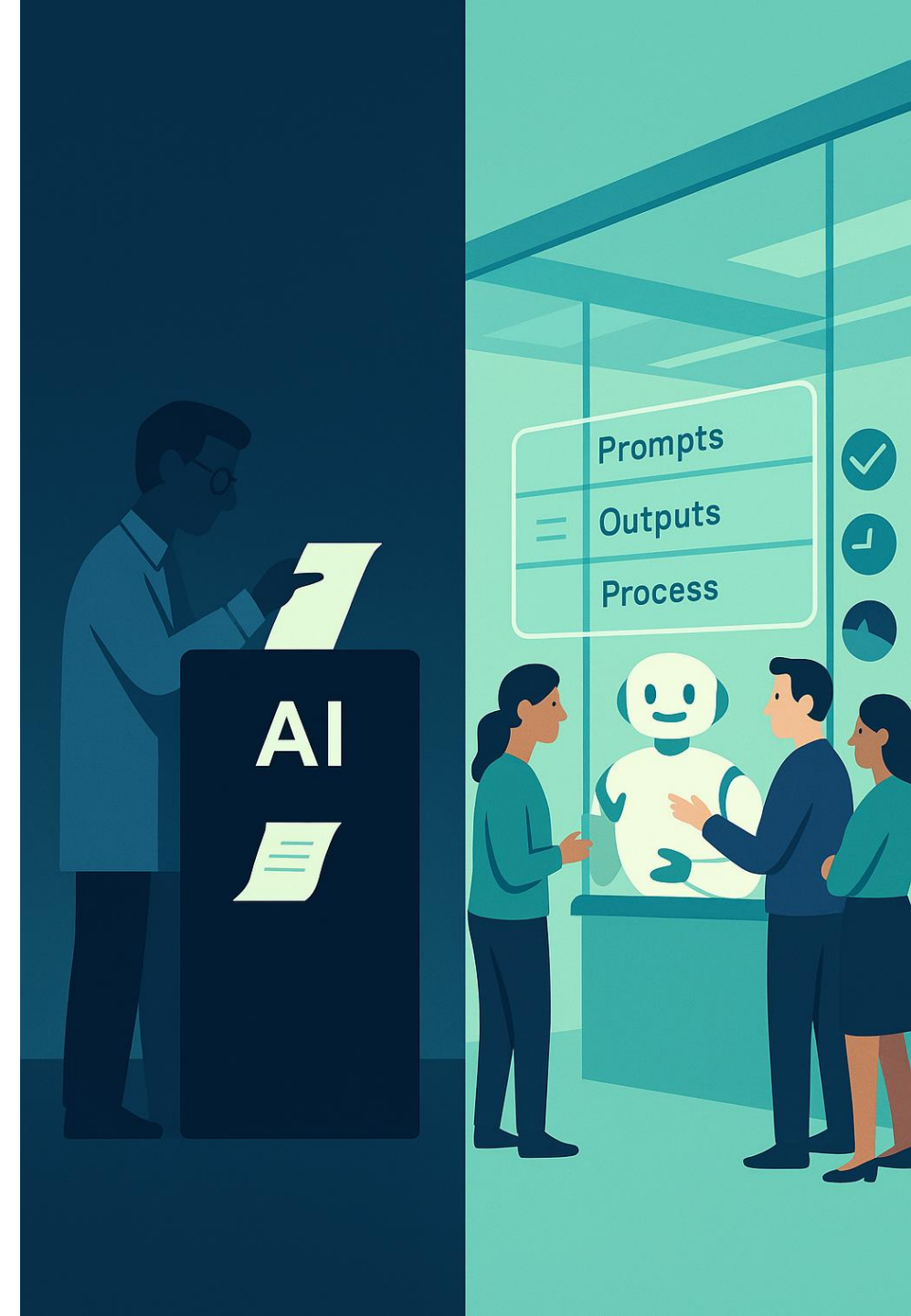
*costlier*

*slower*

# 2. THE DOCUMENTATION

# 1. DOCUMENTAL TRUTH CANNOT/SHALL NOT BE PROVABLE

Non-gameable documentation would be a straightjacket

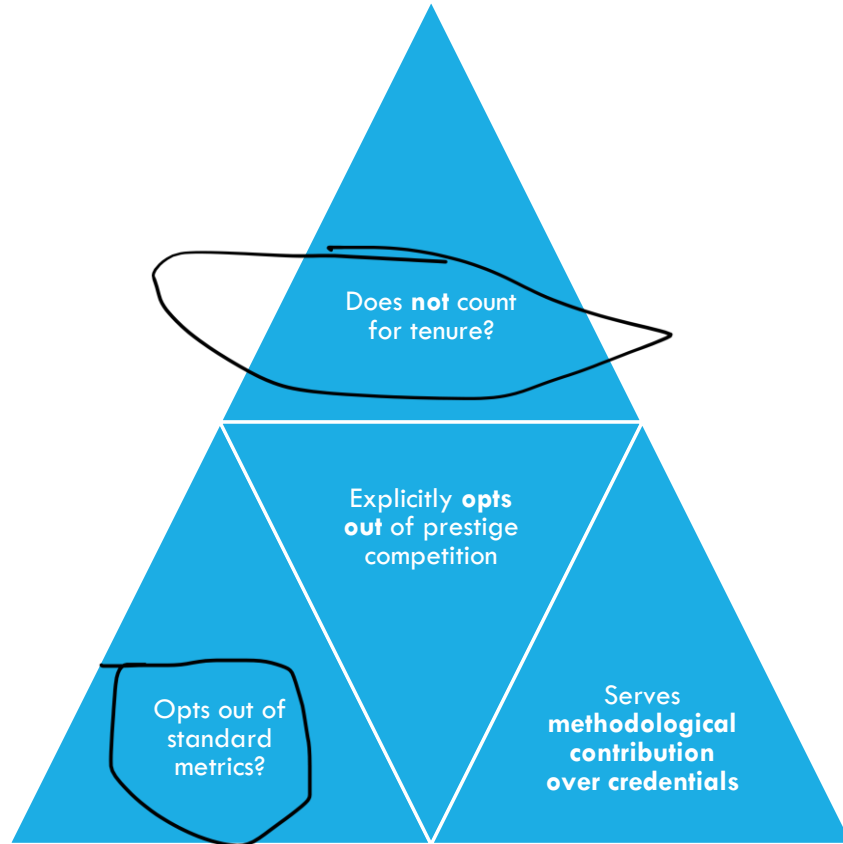Turns philosophers into bureaucrats

Places unreasonable burdens on writing

Non-gameable transparency=creativity killer

# THIS CAN BE ADDRESSED BY REMOVING INCENTIVES TO FAKE TRANSPARENCY (INITIALLY)

# CORE METHODOLOGICAL VALUES



Ecological validity
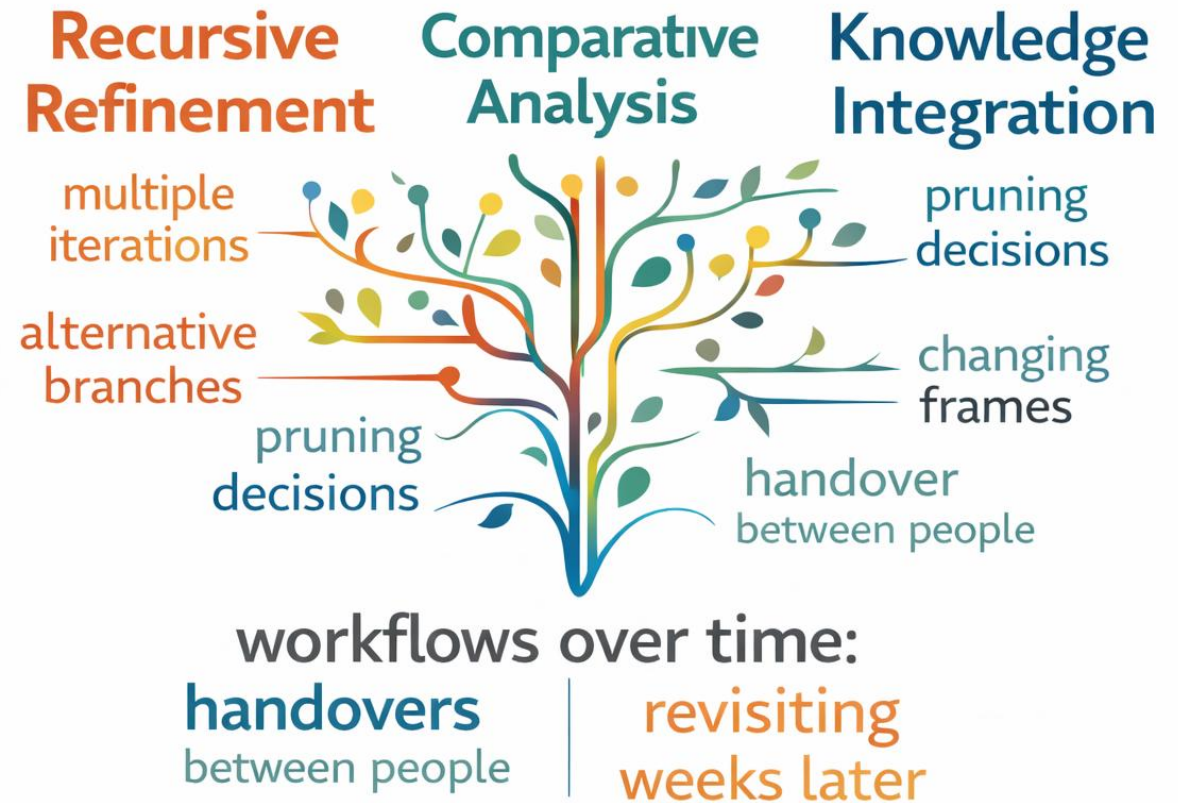


Good faith orientation
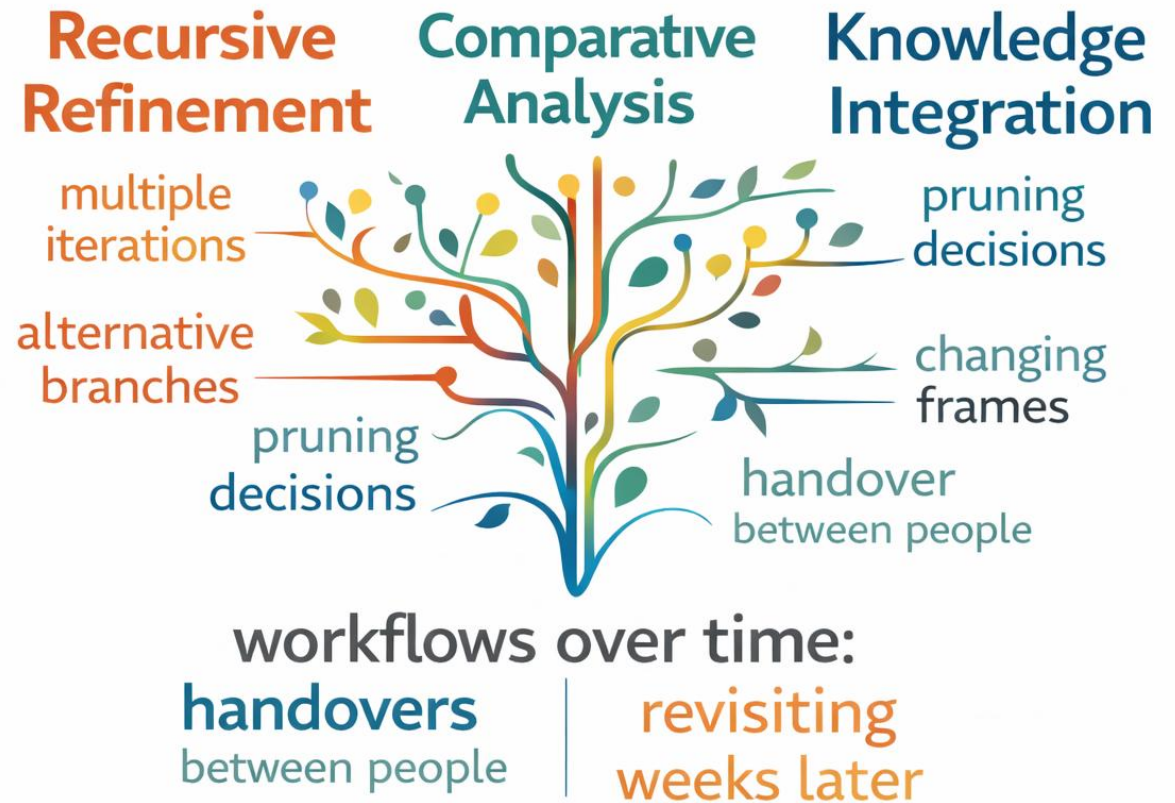
# 2. HONESTY AS THE "PATH OF LEAST RESISTANCE"

Most sophisticated AI-assisted processes are non-linear:
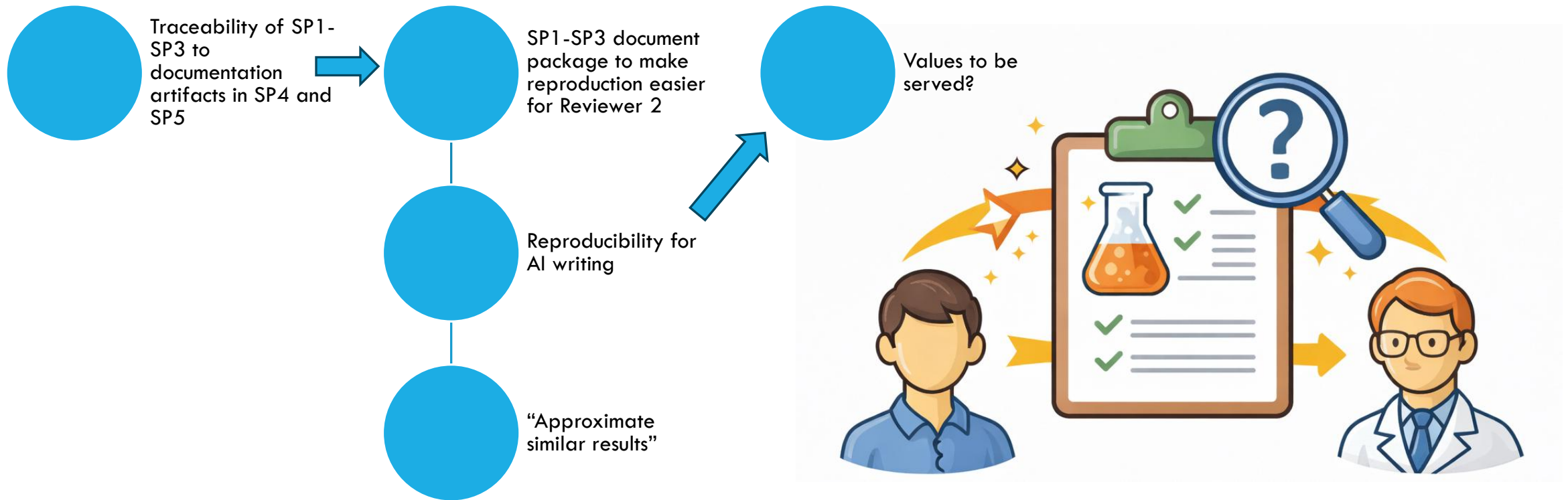
# 2. HONESTY AS THE "PATH OF LEAST RESISTANCE"

Documenting THIS isn't simple.

Faking a complex process (realistically) can be costly.

# A NEW ROLE FOR REVIEWER #2: REPRODUCIBILITY (?)

Traceability of SP1-SP3 to documentation artifacts in SP4 and SP5

SP1-SP3 document package to make reproduction easier for Reviewer 2

Reproducibility for AI writing

"Approximate similar results"

Values to be served?

# A NEW ROLE FOR REVIEWER #2: REPRODUCIBILITY (?)

Current question: is reproducibility essential?

*Is complex/sophisticated + coherent documentation* enough?

Will efficient machine-guided review mechanism for transparency checking be developed?
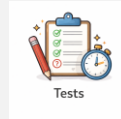
# TOWARDS A STANDARD?

- We need to discuss goals

- We need to discuss standards people are happy to use

- We need run tests

- This can help many people outside philosophy

JPEP APPENDIX

# CONCLUSION

I enjoy conversations with AI about philosophy

These conversations lead to a lot of interesting ideas, some paper worthy

The role of AI in shaping our philosophical ideas should be tracked

Authors' ability to write prompts and shape contexts to deliver good quality prose should become publicly teachable assets

Transparency is needed for both responsibility over ideas (and biases) and learning

Transparency is costly. This can be viewed as a feature, not a bug

Background (unargued) assumption:

      our current peer review processes are running bankrupt

      isolating ourselves from AI progress is foolish

# THANK YOU

This presentation was intentionally unconventional.

In the age of AI, you can produce a traditional synthesis of my paper in slide form with a single prompt

I have chosen, instead, to present the paper to you in a different optic

I wrote this personally, despite my communication style being less effective and more chaotic than ChatGPT's

ChatGPT's generated images have been used