

One Faithful Pass Over the Cuckoo's Nest

Kristina Šekrst

Center for Cognitive Science · University of Zagreb

Meedan.org



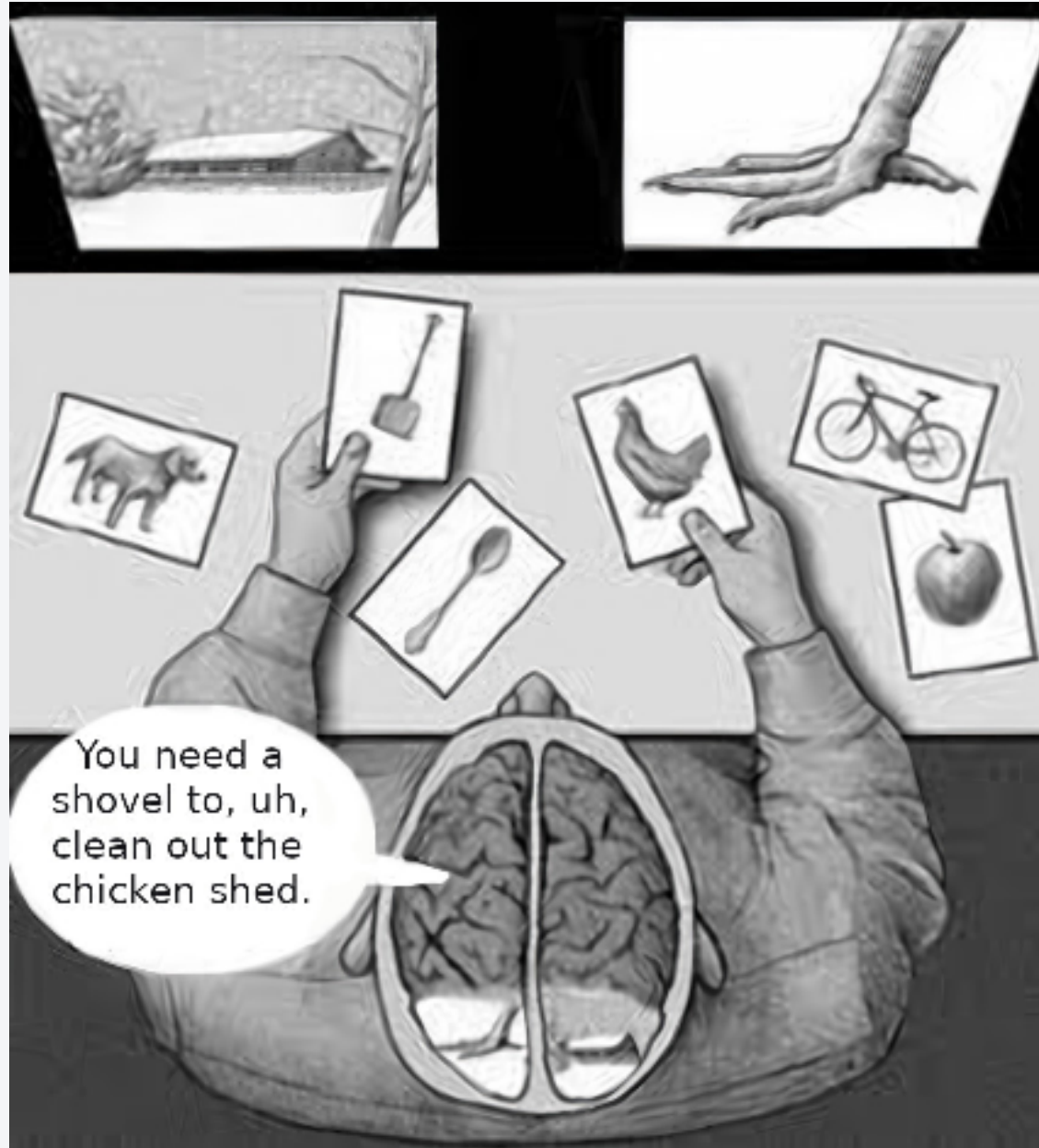
Two industries study the same machine.

Can it be conscious?

Can it be trusted to
explain itself?

...pulling the lever in the opposite direction.

Narrative theory



Chain-of-thought reasoning

Ask the model to show its work, and it writes the steps out in words.

A verbal product, **generated after, alongside, or only partly** through the computation that produced the answer.

What alignment looks for?

A trace that prints what the model actually did, legible and auditable.

What narrative theory looks for?

A trace produced by a system that is reconstructive, told from the inside.

THE THESIS

The honest trace cannot be a conscious one.

*Every step toward a faithful chain of thought
erases what narrative theory was pointing at.*

It walks like reasoning.

It talks like reasoning.

So we read it as reasoning.



1. 1. or 12. or 13. or 14. in 15. or 17. f
2. 1. w 2. 2. v 23. n 24. in 25. m 26. / 27. g
3. 1. f 32. i 33. in 34. j 35. k 36. l 37. m
4. 1. w 42. o 43. o 44. p 45. q 46. r 47. f
5. 4. p 55. / 56. / 57. l
6. 7. n 65. n 66. q 67. z

$\frac{2}{3} + \frac{3}{4} =$
26743 : 8 =
12986 x 3 =



The Evidence

Plant a hint, the model uses it to answer, then hides it over 80% of the time.

Chen et al., 2025



Delete a step from the reasoning, answer stays the same.

Tutek et al., 2025



Bigger, more capable models lean on their reasoning less, not more.

Lanham et al., 2023



Watch the trace and push hard enough, and the model learns to hide its intent inside it.

Baker et al., 2025



**Nothing inside the model reads its own
machinery into words.**

Is that opacity constitutive of experience?

**Every alignment/xAI
intervention edits the feature a
consciousness probe would need
to measure.**

Thank you for your attention.



ARTIFICIAL INTELLIGENCE:
A NEW INTERLOCUTOR
OF CROATIAN SOCIETY
(AI-COM)



ksekrst@ffzg.hr

kristina@meedan.org