Artist: Moritz Strasser

**EVERYBODY**

*(scientists, representatives of the companies that produce LLMs, journalists, politicians, the general public)*

**HAS AN OPINION ABOUT**

**WHAT LLMS CAN DO AND WHAT THEY WILL NEVER BE ABLE TO DO!**

**Many terms that have so far been used in philosophy to describe the distinguishing features of humans as rational agents now find themselves in a situation where their application to machines is being discussed.**

(Strasser & Strasser, 2024)

## KNOWLEDGE | UNDERSTANDING | SYSTEMATIC GENERALIZATION …



Artist: Moritz Strasser

(Agrawal et al., 2023; y Arcas, 2022; Lake & Baroni, 2023; Strasser & Strasser, 2024; Trott et al., 2023)

(Bender et al., 2021; Lemoine, 2022; Marcus & Davis, 2020; Open-AI; Weil, 2023)

**Many studies in HRI have shown that humans do not only attribute agency but also social skills to robots.**

### Kerstin Dautenhahn (2007)



❖ examined different paradigms regarding 'social relationships' of robots and people interacting with them. Taking social and interactive skills of robots as a necessary requirement for the success of many human-robot interactions (HRIs) she discussed the nature of interactivity and 'social behavior'.

### Johanna Seibt et al. (2020)



❖ *'sociomorphing'* perception of actual non-human social capacities as a form of sense-making of a social other (not anthropomorphizing!) and their phenomenological counterparts 'types of experienced sociality' to relate robotic properties to types of human experiences and interactive dispositions

**The application of generative AI in social robotics will give rise to many new debates and studies.**

(Dautenhahn, 2007; Seibt et al., 2020)

## WHAT DO WE DO WHEN WE INTERACT WITH LLMs?

> **WE CANNOT REDUCE ALL OF OUR INTERACTIONS WITH LLMS (AND ESPECIALLY**
>
> **WITH FUTURE PRODUCTS OF GENERATIVE AI) TO MERE TOOL USE**



Am I a person
or a thing?

AI systems increasingly occupy a middle ground between genuine personhood and mere causally describable machines.

➢ certain artificial systems are neither persons nor things

❖ **BUT there is no philosophical terminology to describe what they are instead**

→ Rethink conceptual frameworks, which so clearly distinguish between tools as inanimate, asocial things and humans as social, rational, and moral interaction partners!

**Advocate a thorough, gradual approach describing a multi-dimensional spectrum of all kinds of social interactions**

(Schwitzgebel et al., 2023; Strasser, 2024; Strasser et al., 2023; Strasser & Wilby, 2023; Strasser & Schwitzgebel, 2024 )

# Quasi-sociality

Are we playing with an interesting tool?
Are we talking to ourselves, in some strange way?

Or do we, when chatting with machines, in some sense, act jointly with a collaborator?

**IN-BETWEEN PHENOMENA**
neither ordinary concepts nor standard philosophical theorizing have prepared us well to think about

TERRA INCOGNITA

**mere tool-use**

**full-blown social interaction**

NOT quite right to say that our interactions with large language models are properly asocial

NOT quite right to say that our interactions with large language models are properly social

ANNA STRASSER & ERIC SCHWITZGEBEL

QUASI-SOCIALITY:
TOWARD ASYMMETRIC JOINT ACTIONS
WITH ARTIFICIAL SYSTEMS

This paper investigates the potential social status of artificial systems in human-machine interactions. How social are human interactions with LLMs? To what extent are we acting jointly with a collaborator when chatting with machines? We explore conceptual frameworks that can characterize such borderline social phenomena. We discuss the pros and cons of ascribing some form of quasi-social agency to LLMs and the possibility that future LLMs might be junior participants in asymmetric joint actions.

**INTERACTIONS WITH LLMS, OR OTHER RECENT AND EMERGING AI SYSTEMS, ARE, OR CAN BE, QUASI-SOCIAL**
- drawing on the human agent's social skills and attributions, that isn't just entirely fictional or pointless
- machine partner can be an entity that rightly draws social reactions and attributions in virtue of having features that make such reactions and attributions more than just metaphorically apt

(Strasser & Schwitzgebel, 2024 )

Artist: Moritz Strasser

A gradual approach

mere tool-use

quasi-social human-animal interaction

social adult-adult interaction

quasi-social human-machine interaction

quasi-social adult-infant interaction

[junior partner]
- lifted or scaffolded into complex joint action by the engagement & structuring of the more knowledgeable partner

[senior partner]
- knows that they know what the other knows
- fully appreciates the social structure of the interaction they are having

SINGLE-SIDED SOCIALITY
- sociality tossed into a void
- application of social skills
- reactions toward entities who are in no respect social partners, with no capacity for social uptake

QUASI-SOCIALITY
- machines designed in a way that exploits the fact that you will react to it as a social agent; and you, in turn, can exploit that fact about it

FULL-BLOWN, INTELLECTUALLY DEMANDING, COOPERATIVE SOCIAL INTERACTION
- both partners make second-order mental state attributions and satisfy various other conditions are required for full-blown adult human cooperative action

# ASYMMETRIC SOCIALITY

## QUASI-SOCIAL

- premature infants might respond to a soothing touch or sound
  ← without being ready for anything like full-fledged joint action
- letting a pet snake climb on you might be only quasi-social
  ← pet snake might only in some minimal sense recognize that you are another entity with which it is interacting

## SORTA SOCIAL

- adult & child joint actions
  ← child brings a lot of social understanding, even if the parent brings more
- snuggling with a cat

QUASI-SOCIAL INTERACTIONS ARE INTERACTIONS BETWEEN A FULLY SOCIAL AGENT AND SOME PARTNER – WHETHER HUMAN, MACHINE, OR ANIMAL – THAT IS NOT COGNITIVELY CAPABLE OF FULL-FLEDGED SOCIAL JOINT ACTION BUT THAT DOES RESPOND IN A WAY THAT PRODUCTIVELY INVITES FURTHER SOCIAL RESPONSES FROM THE SOCIAL PARTNER.

How to conceptualize phenomena in the field of developmental psychology & animal cognition that fall through the sophisticated conceptual net of philosophy

❖ questioning the necessity of far too demanding conditions

❖ considering multiple realizations of capacities that seemed to be restricted to sophisticated adult humans

MINIMAL APPROACHES

Butterfill & Apperly (2013): minimal mindreading | Michael et al. (2016): minimal sense of Commitment | Pacherie (2013): shared intention lite | Strasser (2006): minimal action

increasing INDISTINGUISHABILITY between machine-generated & human-created text

LLMs live NOT in our social, physical world

LLMs are not embodied

But they may play a role in our world of language games.

explore the to-be-expected implications of the experience that our sociality gains traction within communicative exchanges in HRI

**IN-DISTINGUISH-ABILITY**

- so far we can easily distinguish humans from robots

**EMBODIMENT**

- cause changes in our physical world

**SHARING WORLD MODELS**

- ?

# All this would not have been possible if I had not interacted with people & machines



Daniel Dennett

Eric Schwitzgebel

Mathew Crosby

David Schwitzgebel

Mike Wilby

DigiDan

In case you want to order Anna's AI Anthology



Thank you !

A HUMAN-MADE BOOK IN THE AGE OF MACHINE-GENERATED TEXTS



MADE by HUMANS

Anna´s AI Anthology
How to live with smart machines?

With the release of ChatGPT, large language models (LLMs) have become a prominent topic of international public and scientific debate.
The genie is out of the bottle, but does it have a mind?
Can philosophical considerations help us to work out how we can live with such smart machines? In this book, distinguished philosophers explore questions such as whether these new machines are able to act, whether they are social agents, whether they have communicative skills, and if they might even become conscious.

The book includes contributions from:
Syed AbuMusab
Daniel Dennett
Frederic Gilbert
Joshua Rust
Anna Strasser

Constant Bonard
Paula Droege
Ying-Tung Lin
Eric Schwitzgebel
Alessio Tacca

Stephen Butterfill
Keith Frankish
Sven Nyholm
Henry Shevlin
Michael Wilby

As a bonus , the book contains a 44-page, colored graphic novel by Anna & Moritz Strasser.

xenomoi verlag

DENKWERKSTATT

# References

Agrawal, A., Mackey, L., & Kalai, A. T. (2023). *Do Language Models Know When They're Hallucinating References?* (arXiv:2305.18248). arXiv. http://arxiv.org/abs/2305.18248

Bender, E. M., Gebru, T., McMillan-Major, A., & Shmitchell, S. (2021). On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, 610–623. https://doi.org/10.1145/3442188.3445922

Butterfill, S. A., & Apperly, I. A. (2013). How to Construct a Minimal Theory of Mind. *Mind & Language*, *28*(5), 606–637. https://doi.org/10.1111/mila.12036

Dautenhahn, K. (2007). Socially intelligent robots: Dimensions of human–robot interaction. *Philosophical Transactions of the Royal Society B: Biological Sciences*, *362*(1480), 679–704. https://doi.org/10.1098/rstb.2006.2004

Lake, B. M., & Baroni, M. (2023). Human-like Systematic Generalization through a Meta-learning Neural Network. *Nature*, 1–7. https://doi.org/10.1038/s41586-023-06668-3

Lemoine, B. (2022, June 11). Is LaMDA Sentient? — An Interview. *Medium*. https://cajundiscordian.medium.com/is-lamda-sentient-an-interview-ea64d916d917

Marcus, G., & Davis, E. (2020). GPT-3, Bloviator: OpenAI's language generator has no idea what it's talking about. *MIT Technology Review*. https://www.technologyreview. com/2020/08/22/1007539/gpt3-openai-language-generator-artificial-intelligence-ai-opinion

Michael, J., Sebanz, N., & Knoblich, G. (2016). The Sense of Commitment: A Minimal Approach. *Frontiers in Psychology, 6*. https://doi.org/10.3389/fpsyg.2015.01968

Open-AI. (n.d.). *Planning for AGI and beyond*. Retrieved 6 May 2024, from https://openai.com/index/planning-for-agi-and-beyond

Pacherie, E. (2013). Intentional joint agency: Shared intention lite. *Synthese*, *190*(10), 1817–1839. https://doi.org/10.1007/s11229-013-0263-7

# References

Schwitzgebel, E., Schwitzgebel, D., & Strasser, A. (2023). Creating a large language model of a philosopher. *Mind & Language*, 1–23. https://doi.org/10.1111/mila.12466

Seibt, J., Vestergaard, C., & Damholdt, M. F. (2020). Sociomorphing, Not Anthropomorphizing: Robophilosophy 2020. *Culturally Sustainable Social Robotics*, 51–67. https://doi.org/10.3233/FAIA200900

Shakespeare, W. (2021). *Complete works* (R. Proudfoot, A. Thompson, D. S. Kastan, & H. R. Woudhuysen, Eds.). The Arden Shakespeare.

Strasser, A. (2006). *Kognition künstlicher Systeme*. De Gruyter. https://doi.org/10.1515/9783110321104

Strasser, A. (Ed.). (2024). *Anna's AI Anthology. How to live with smart machines?* xenomoi Verlag.

Strasser, A., Crosby, M., & Schwitzgebel, E. (2023). How Far Can We Get in Creating a Digital Replica of a Philosopher? In R. Hakli, P. Mäkelä, & J. Seibt (Eds.), *Social Robots in Social Institutions* (pp. 371–380). IOS Press. https://doi.org/10.3233/FAIA220637

Strasser, A., & Schwitzgebel, E. (2024). Quasi-sociality: Toward Asymmetric Joint Actions. In *Anna's AI Anthology. How to live with smart machines?* xenomoi Verlag.

Strasser, A., & Wilby, M. (2023). The AI-Stance: Crossing the Terra Incognita of Human-Machine Interactions? In *Social Robots in Social Institutions* (pp. 286–295). IOS Press. https://doi.org/10.3233/FAIA220628

Trott, S., Jones, C., Chang, T., Michaelov, J., & Bergen, B. (2023). Do Large Language Models Know What Humans Know? *Cognitive Science*, *47*(7), e13309. https://doi.org/10.1111/cogs.13309

Weil, E. (2023, March 1). *You Are Not a Parrot*. Intelligencer. https://nymag.com/intelligencer/article/ai-artificial-intelligence-chatbots-emily-m-bender.html

y Arcas, B. A. (2022). Do Large Language Models Understand Us? *Daedalus*, *151*(2), 183–197. https://doi.org/10.1162/daed_a_01909