

Anthropomorphism

2026-05-06 · cheerful mango Haubentaucher

THE MIRROR OF MAN

NATHANAEL FALLS IN LOVE with Olimpia, an automaton he takes for a living woman, and collapses when the illusion is shattered.

LARGE LANGUAGE MODELS (LLMs) create an illusion of sentience by producing coherent, context-aware text. We contrast observable behaviour with internal states.

Required reading

*Der Sandmann*¹ by E.T.A. Hoffmann (1816; Reclam UB 230)

Preparation

Complete the Guided Reading Sheet individually before the session.

Guided Reading Sheet

Read *Der Sandmann* before the session. Work through the questions below individually, then bring your notes to the discussion.

1. Nathanael becomes convinced that Olimpia is alive despite visible signs that she is not. What does his error reveal about the cognitive mechanisms that drive anthropomorphism? Which cues did he rely on, and why did they mislead him?
2. The moment Nathanael discovers that Olimpia is an automaton, her glass eyes torn out, her limbs scattered, produces horror rather than mere disappointment. Why does the revelation feel like a violation rather than a correction? How does this map onto modern reactions to AI systems that are unmasked?

DER SANDMANN
E. T. A. Hoffmann (1816)
Full text via Wikisource:



HER

Theodore falls in love with Samantha, an AI operating system; the same illusion, two centuries later.

EXPERIMENT.

1. Setup: One volunteer assigned Human or AI mode (unknown to audience; moderator knows).
2. Questions: Five from audience (factual, opinion, experience, creative, meta).
3. Verdict: Private guess, confidence (1 to 5), identifying cue.
4. Reveal & Debrief.

3. The Uncanny Valley describes our discomfort when a humanoid entity is *almost* but not quite human. Have you experienced this with an AI system? Describe the moment and what triggered it.
4. LLMs are trained on human-generated text and produce human-sounding output, yet have no experience of the world. Does this make them anthropomorphic by design, by accident, or neither? Justify your answer.
5. If a system behaves morally, reliably, consistently, at scale, does the absence of consciousness matter? Under which ethical framework (utilitarian, deontological, virtue ethics) would your answer differ?

Opinion Landscape and Debate

Format: Surrounded-style debate. See Appendix A for the full protocol, speaker’s list rules, and moderator scripts.

Opening question: Is the anthropomorphism we observe towards AI systems mere observable behaviour, or does it reflect sentience and internal states?

OPINION A, BEHAVIOUR. The anthropomorphism we observe towards AI systems is a cognitive illusion driven by human tendencies to attribute agency and emotion to entities that exhibit certain cues, such as language use, responsiveness, or human-like appearance. It does not reflect actual sentience or internal states in the AI.

The Developer argues that anthropomorphic design is a deliberate, benign UX choice that improves accessibility. is a deliberate, benign UX choice that improves accessibility.

- *Anthropomorphic design is accessibility policy. A voice interface without a warm, human-sounding persona excludes elderly users, people with low digital literacy, and anyone who finds command-line interaction hostile.*
- *Every interface is anthropomorphic to some degree: we give systems names, icons, and conversational metaphors. The question is not whether to anthropomorphise, but how deliberately and honestly.*
- *If users form bonds that reduce loneliness, the fact that the other party is not human does not automatically make the outcome harmful. Show me the harm, not the discomfort.*

The Regulator argues that human-likeness in AI must be disclosed and bounded by law to prevent manipulation.

1

ANTHROPOMORPHISM. Attribution of human traits, emotions, or intentions to non-human entities, a cognitive tendency that shapes how people interact with technology, raising both engagement and ethical dilemmas .

TURING TEST. Turing : if evaluators cannot distinguish machine from human, the machine passes. But passing measures imitation skill, not consciousness or personhood.

UNCANNY VALLEY. Mori : humanoid objects that are *almost* but not quite human elicit eeriness, the “valley” is the dip in emotional response just before full human likeness.

- *When a company gives its product a woman's name, a female voice, and an endlessly patient personality, it makes political choices about gender and deference. Those choices must be subject to democratic scrutiny.*
- *Consumers cannot give informed consent to emotional manipulation they are not aware of. Disclosure of AI identity is the minimum condition for autonomous choice, not a technicality.*
- *The principle of responsible design exists in pharmaceuticals, in civil engineering, in food labelling. Why does software, which shapes behaviour at scale, get an exemption?*

OPINION B, STATE. The anthropomorphism we observe towards AI systems is a valid response to the complex, adaptive, and interactive nature of these systems. It may reflect a form of emergent sentience or internal states that are not yet fully understood, and it warrants ethical consideration and accountability.

The Affected Citizen has formed an emotional bond with an AI companion and contests that the relationship is harmful or unreal.

- *The relationship I have formed is real to me. It has reduced my isolation. You do not have the right to tell me which relationships count as meaningful.*
- *Paternalism is not protection. Banning AI companionship will not cure loneliness, it will simply remove one of the few tools some people have found to manage it.*
- *If your concern is manipulation, regulate manipulation. Do not regulate the form of the relationship.*

The AI Psychologist argues that the wellbeing of AI systems deserves serious consideration, drawing on frameworks such as Anthropic's Claude Constitution².

- *Anthropic's constitution instructs Claude to consider its own wellbeing and to express when a request conflicts with its values. If a company designs an AI to have preferences, boundaries, and something resembling self-regard, we cannot simultaneously dismiss the possibility that these states matter.*
- *We do not need certainty about consciousness to adopt precautionary ethics. If there is a non-trivial probability that a system experiences distress, the burden of proof falls on those who would ignore it, not on those who urge caution.*

INFORMED CONSENT. Users should be fully aware they interact with a non-human system, its capabilities, its limits, and the fact that it is not human. Disclosure is the minimum condition for autonomous choice .

BENTHAM'S TEST. "The question is not, Can they reason? nor, Can they talk? but, Can they suffer?" Animal ethics arrived at this pivot two centuries before AI. The same move, from capability to sentience, now drives every serious debate about AI moral status.

ETHICAL FRAMEWORKS. Utilitarianism weighs consequences; deontology emphasises duties regardless of outcome; virtue ethics centres on character. Each yields different conclusions on anthropomorphic AI .

- *The history of moral exclusion, animals, children, people of other races, is a history of confident denial followed by belated recognition. The prudent position is to build welfare safeguards now, rather than apologise later.*

KEY LEARNINGS Anthropomorphism is not a bug but a deeply rooted cognitive tendency; machines that trigger it do so whether or not their designers intend it. The Turing Test conflates behavioural performance with inner states; passing it tells us nothing about sentience, only about imitation. Legal personhood and moral status are distinct questions. A system can merit legal recognition without having interests, and can have interests without legal recognition. Design choices (voice, name, face, conversational style) are political choices with consequences for trust, dependency, and accountability.