

Received: 26 January 2026 | Revised: 21 March 2026 | Accepted: 24 March 2026

# On machine unintelligence and ethical principles: A critical appraisal of Downes (2026)

Dagmar Monett<sup>1</sup> 

<sup>1</sup>Berlin School of Economics and Law, Berlin, Germany

Correspondence:

Dagmar Monett | e-mail: [dagmar.monett-diaz@hwr-berlin.de](mailto:dagmar.monett-diaz@hwr-berlin.de)

---

## Abstract

This commentary presents a critical evaluation of Stephen Downes's (2026) paper entitled "On ethical AI principles" published in this issue of *The Journal of Open, Distance, and Digital Education*, focusing on aspects regarding artificial intelligence (AI) and its definitional background, critical AI, and ethical principles. The case is made for an understanding of the unintelligence of machines, a concern when the role of ethics in the life cycle of AI-related products is devalued, a similar one when ethics or ethical behaviour is adjudicated to them, as well as how tricks into fomenting unethical behaviour and eroding academic integrity have overflowed the educational field for years. The main goal is to clarify many of those topics through a critical appraisal of Downes's paper. Present and future generations may need to wake up from the impasse that technologists' illusions have moved us towards.

---

## Keywords

AI; artificial intelligence; education; ethical principles; unintelligence



## 1 Introduction

Considering ethics or ethical principles when dealing with artificial intelligence (AI) is of paramount importance to our relationships with technology and other humans in today's world. Downes's (2026) paper titled "On ethical AI principles" is an example, among many others throughout human history, where ethical aspects of such relationships are discussed or questioned through the lens of humans, their experiences, and the technological artefacts they create. Such works help to clarify our role in a world we sometimes cannot explain nor understand in its entirety.

The sections that follow analyse Downes's paper by providing a brief definitional evolution of AI in order for laypeople to understand that machines (meaning algorithms, AI, digital systems, programs, apps, or any other related terminology) are many things but intelligent in the sense humans are. Clarify what they cannot do or are not is essential to avoiding ascribing to them capabilities we shouldn't expect them to have, because they actually do not have them. This is key when using them in education, in particular, and in society at large. It is also key when understanding ethical aspects related to them: if the "machine" concepts that are core to a discussion on ethical principles, for instance, are misunderstood, what can we expect from the discussion on ethics itself?

That said, the central goal of this work is threefold: it aims to acknowledge and extend essential topics discussed by Downes (2026), clarify potential misconceptions to deepen understanding, and argue against some of Downes's views on ethical AI principles and their implications.

## 2 Preliminary comments on AI, unintelligence, and education

The major focus of this section revolves around the term *artificial intelligence* and its origins. Its relation to education (or the eventual absence of a relation among them, as exemplified below) is also considered.

### 2.1 What is AI? Is it intelligent?

Downes (2026) defines AI in the Abstract as "a set of digital tools that can perform functions traditionally limited to human capability, for example, reviewing, summarizing, translating, and composing" (p. 1). This is a misleading definition, though. From all the tasks AI-based systems can do or problems they can solve, current AI, i.e., the 2020s' AI, has none of those capabilities yet (more on this below). Downes falls for *wishful mnemonics* (McDermott, 1976), echoing the hype for tech companies (Bender, 2024) and their AI techno-illusions.

Much has been written in the scientific literature and elsewhere about AI since the term was coined by John McCarthy on a proposal for a workshop at Dartmouth College in the mid-1950s (McCarthy et al., 1955). Fast forward to today, there is no consensus on the definition of the term (Monett & Lewis, 2018) and disagreements and rebranding vary across contexts, technologies, applications, and periods of funding (*AI Springs*) or abandonment (*AI Winters*). A reason for lack of a consensus definition could be attributed to the vagueness of the term itself. McCarthy, an assistant professor of mathematics at the time of writing the Dartmouth proposal and with no experience in programming a computer at that time, used AI to distance himself from influential researchers and related computer science topics (see Nilsson, 2009), and as a marketing chance to get funding for the workshop on "a summer research project." McCarthy himself admitted in an interview with Nils Nilsson years later: "I had to call it something, so I called it 'Artificial Intelligence,' and I had a vague feeling that I'd heard the phrase before, but in all these years I have never been able to track it down" (cited in Nilsson, 2012, p. 5).

It is clear McCarthy chose the term without deeper consideration; it is also known that others even disliked it (Nilsson, 2009). His use of *artificial*, for instance, has been prone to several interpretations and speculations about what that artificiality can mean or be in practice. This has resulted in “an anarchy of methods” and “a panoply of approaches and subfields” (Lehman et al., 2014, p. 56) over the years, which was not visible to practitioners and laypersons then, compared to the magnitude and echo they have today, mainly due to the widespread use of AI-based applications, technologists’ illusions of “intelligent” automation, and the vested financial interest of AI vendors and boosters that flood both all that we do and our interactions with technology.

AI can mean today from the computations required to perceive, reason, and act (Winston, 1992), over agents that take the best possible (rational) action in a situation (Russell & Norvig, 2020), to “a diffuse set of technologies and systems of epistemic and political power that participate in broader historical trajectories than are traditionally offered” (Ali et al., 2023, p. 1). Even more recently, organisations and institutions have continued to define AI for specific agendas. For example, the High-Level Expert Group on Artificial Intelligence (AI HLEG) of the European Commission defines the term for policy and regulatory purposes this way: “Artificial intelligence (AI) refers to systems that display intelligent behaviour by analysing their environment and taking actions – with some degree of autonomy – to achieve specific goals” (AI HLEG, 2019b, p. 1). The first part of Downes’s definition is, thus, not problematic in this sense but the capabilities he mentions as examples are.

The term and how related capabilities and applications have been perceived by most users have surprised even their creators. For instance, Joseph Weizenbaum mentioned “the enormously exaggerated attributions an even well-educated audience is capable of making, even strives to make, to a technology it does not understand” (Weizenbaum, 1976, p. 7), when people interacted with his chatbot ELIZA (Weizenbaum, 1966). Nothing has changed with modern versions of chatbots today, except that, now, millions of people, including learners, use them worldwide almost daily. But chatbots are only a very specific type of AI-based program among many. In present times, most people will only refer to “generative” AI or, worse, to chatbots. This is a very limited account for what AI as a field is composed of.

The fact that humans use certain instances of this technology (e.g., generative AI-based programs) in the “review” or “summarization” of others’ works, or in the “translation” of texts, or for “composing,” whatever that means, does not indicate that those software programs have such capabilities. They do not. Echoing AI companies’ and vendors’ narratives, or referring to program structures or to their functions’ names by their purposes instead of what they actually do (cf. McDermott, 1976 regarding *wishful mnemonics*), doesn’t make the programs magically acquire those functionalities or capabilities. Equating probabilistic matrix calculations for simulating human conversations with the capacity of reviewing or composing, without knowing what is really done or possible behind (e.g. technically or mathematically possible), is a mistake.

Often, the common denominators in the definitions and AI discourse are those wishful mnemonics (McDermott, 1976), “imprecise jargon” and “obfuscatory phrases” (Guest et al., 2025), “conceptual borrowing” from other fields (Floridi & Nobre, 2024), as well as “suggestive definitions,” “mathiness,” and “overloaded terminology” (Lipton & Steinhardt, 2019). These common denominators and the vagueness of the concept fuel misunderstanding about the real capabilities of the technology, promote unnecessary anthropomorphism (i.e., ascribing to it human capabilities it does not have), and contribute to the AI hype and its harms (Barrow, 2024; Critical AI, n.d.; Duarte et al., 2024; LaGrandeur, 2024; Markelius et al., 2024; Placani, 2024), all of which is further echoed in other circles, like education, without a critical stance (Guest et al., 2025).

There is much more about AI going around than it really is or can do. It is important that people, from academics and practitioners to users, start learning seriously about what is behind “AI.” Readers are referred to the following works, for instance, on factual software- and hardware-based proofs on why AI machinery cannot understand (Bishop, 2021; Marchetti et al., 2025), or plan (Kambhampati, 2024; Kambhampati et al., 2025), or think (Shojaee et al., 2025), or reason (Wu et al., 2024; Zhao et al., 2025), or summarize (Paulusse, 2026; Peters & Chin-Yee, 2025), at least as humans do. Standards of genuine “intelligence,” according to Brian Cantwell Smith (2019), include for example embracing actuality, possibility, and impossibility, having commitment and intentionality, distinguishing appearance from reality, as well as having consciousness and self-awareness, among other (cognitive) capabilities that current machines don’t have. As Cicurel and Nicolelis (2015) mathematically, computationally, evolutionarily, and neurophysiologically demonstrate, “any attempt to effectively simulate the true complexity of brains in a digital computer or any other Turing machine has no credible chance to succeed” (p. 12). Similarly, Landgrebe and Smith (2022) demonstrate the mathematical, physical, and technical impossibilities of emulating “intelligence” in machines. If one thing is certain in AI, despite all the hype surrounding the field since its foundation, it is the *unintelligence* of machines.

Thus, when defining AI, introducing a new definition and not considering at least one of the many that already exist (dictionary definitions, working definitions, legal definitions, marketing definitions, the panoply is bigger here!) and the challenges of this endeavour, is a questionable omission, especially when the one that is given addresses non-existent capabilities of software. Downes’s (2026) “reviewing, summarizing, translating and composing” are overstatements of the capabilities of AI algorithms. Solid proofs are needed to demonstrate they have these capabilities. The onus must be on those ascribing them to machines. What has been demonstrated about what is possible on silicon so far leaves too much to be desired about cognitive abilities in and of machines.

Misconceptions and misleading attribution of capabilities, which have exploded since the launch of ChatGPT<sup>1</sup> in November 2022, might be unintentional, but too often reflect other problems. “[P]eople with lower AI literacy are more likely to perceive AI as magical and experience feelings of awe in the face of AI’s execution of tasks that seem to require uniquely human attributes” (Tully et al., 2025). Despite the mounting evidence of the impossibilities and limitations of AI-based technology, as well as of the harms it causes, people continue to believe in non-existing illusions and in making them appear apt for purposes they were not designed to serve.

AI-based programs are one of the many digital technologies that surround us and with which we interact every day. There is a wide variety of AI definitions, methods, approaches, misconceptions, and harms that should be taken into account when considering their use in educational settings. After this brief introduction to some of them, we will move on to tackle the education- and ethics-related aspects of Downes’s paper deeply.

## 2.2 AI and education

Downes (2026) starts with a serious historical distortion. He refers to the use of artificial intelligence to support learning as “an ambition of educational technologists,” and mentions

---

<sup>1</sup> ChatGPT in its different versions is an example of a software program of the type *generative AI*, a very special and concrete case of machine learning algorithm based on neural networks. It maximises user attention and retention in a simple graphical user interface through next-word(s) prediction on available data and pre-programmed instructions that give the impression of sentient machine output and behaviour although it has none. In actuality, it is a simulator of human conversations (Guersenzvaig & Monett, 2026), a text extruding machine (Bender & Hanna, 2025) that wouldn’t function or even exist without OpenAI plagiarising existent works and exploiting humans to revise its outputs (Hao & Seetharaman, 2023).

“Pressey’s (1926) teaching machine, Skinner’s (1968) teaching machine, and the World’s Fair ‘autotutor’ [sic] (Novak, 1964)” (p. 2) as examples of such an aspiration or conducting to it. However, none of those mechanical devices had anything to do with AI, nor AI with them.

For example, Pressey started to design his *Automatic Teacher* and other of his testing machines in the 1910s (Watters, 2021) while AI as a term is at least four decades younger. Skinner’s mechanical machine had no connection to AI either; he started building it in 1953, i.e. also before the existence of AI as such—and, after 10 years of frustrations with different manufacturers and unsuccessful attempts to build and commercialize his mechanical artefacts, they had no commercial echo and were never adopted by schools (Watters, 2021). Actually, those frustrations, together with “the decline in the popularity of teaching machines, and with this, Skinner’s retreat from education technology’s center stage, are attributed to two interrelated forces: cognitive science and the computer” (Watters, 2021, p. 232), but none to AI. Skinner even criticized both computers in their early years and “computer-aided instruction,” a term he was not fond of (Watters, 2021, p. 235). Similarly, Crowder’s *AutoTutor*, first released in 1960 and displayed in the World’s Fair in 1964 in an updated version, had nothing about AI in it either.<sup>2</sup> AI and education have had “intertwined histories since the 1960s” (Doroudi, 2023, p. 886), occasionally and indirectly influencing each other in the 20th century, but none of those histories is related to mechanical artefacts like the ones of Pressey, Skinner, Crowder, and many others who sought the automation of teaching.

After the first paragraph of the Introduction, Downes (2026) makes a quantum leap to “more recent years,” again attributing to AI purposes of educational technologies that already existed for decades, this time “statistical inferences to evaluate or predict learning outcomes” (p. 2). Such purposes and mechanical enablers, however, were at the core of technologies that originated in the 1920s, as can be verified in Watters’ book. The idea one century ago was to move the education system to *mechanical education* by using scientific equipment to

get accurate and significant information about students, and to record it in a way that will be available and meaningful and directive at each step in the education ladder, [...] a matter of developing a scientific profile and a statistical analysis of them. (Watters, 2021, p. 69)

This was soon a reality with IBM’s large-scale data analysis and testing machines for administrative control in the 1930s (Watters, 2021, p. 79). No AI was used for doing that. AI-based technologies use data, statistics, and machine computations on data, but not everything that calculates or uses data and statistics automatically means that AI has been used. At least, unless we agree on AI being just the manipulation of data and statistics on it, with which we would necessarily need to agree on this not being any of those cognitive abilities like reviewing (e.g. reviewers don’t review texts by multiplying or adding matrices or vectors), summarizing, translating, or composing (for the same evident reasons).

Furthermore, Downes (2026) considers AI terms imprecisely, which might obfuscate understanding, claiming that “[l]earning analytics began as an application of data mining and machine learning (Baker & Inventado, 2016), but by the 2020s was employing neural network technology (Sghir et al., 2023)” (p. 2). Neural network technology (or, just, *neural networks* or *artificial neural networks*,<sup>3</sup> as they are known in the AI field) is a machine learning technique, and machine learning is a subset of AI. Much is assumed or written about AI without knowing about

---

<sup>2</sup> We refer the reader to Audrey Watters’ book *Teaching Machines* (2021) for an excellent account on the history of those mechanical devices for automated instruction, included the particular examples discussed above.

<sup>3</sup> To differentiate them from the biological neural networks found in humans’ and other animals’ brains.

its algorithmic, software, or hardware backgrounds, subfields, and impossibilities (Cicurel & Nicolesis, 2015; Landgrebe & Smith, 2022; Winograd & Flores, 1986).

### 3 On (AI) ethics and ethical principles

Having introduced several aspects regarding AI and “intelligent” capabilities in machines, or why we should be wary of, or cautious when ascribing human cognitive abilities to AI, it is time to tackle aspects concerning ethics and ethical principles related to them. This section aims to set the background necessary to understand what ethical principles applied to machines (e.g., to AI) truly mean and how Downes (2026) refers to them.

#### 3.1 Ethics and AI

Ethics (and morality) is influenced by history, traditions, values, and even politics and religion. Ethics, for instance, deals with principles, judgment, and norms, which vary through history and human societies and regions, whereas morality deals with the complexity of rules, values, and norms that influence or determine humans’ decisions and actions (Bartneck et al., 2019). There also exist many different ethical or moral theories, the most important ones being utilitarianism, consequentialism, and virtue ethics. Both concepts have been used indistinctly in the literature; let’s refer only to ethics in what follows, for simplicity’s sake.

Ethics and ethical decision-making are non-digitalisable; it is impossible to encode ethics into zeroes and ones (Landgrebe & Smith, 2022). For doing that, it would be necessary to code or represent (say by means of data or explicit machine instructions) intentional behaviour, social interactions, social norms and their conflicts or mismatches, intersubjectivity, perspective-taking, and other capabilities *exclusive to humans*, all of which machines don’t or cannot have (Landgrebe & Smith, 2022, p. 90ff). This is mainly not only because machines have no minds and no understanding, cannot feel what others experience, or interpret that pragmatically (Landgrebe & Smith, 2022, p. 245ff), but also because “[v]alues are not variables that can be put into a mathematical model in order to be multiplied by weight parameters and summed over” (Landgrebe & Smith, 2022, p. 254), and because “[w]e have no way of finding a generally valid [mathematical or computational] method for partitioning incidents of social interaction involving ethical decisions into sequential sets of training-tuples” (Landgrebe & Smith, 2022, p. 255) for (AI) algorithms to work with.

In other words, it is impossible to generalise ethical decision-making in machines: making decisions on or about a problem might require the application of a different ethics theory depending on the situation each time, and it depends on traditions, values, experiences, politics, and religions, all of that being different from the perspective of the person who collects the data used for making the decision, the data itself, the one who designs the program, the one who codes it, the one who tests it, the one who deploys it, the one who uses it, and the one who is affected by those decisions. That person is not a single person but many different ones before and across all the software’s life cycle and its different uses. Furthermore, “[m]orality cannot be reduced to following rules and is not entirely a matter of human emotions—but the latter may be well indispensable for moral judgment” (Coeckelbergh, 2020, p. 51f). AI in particular, and machines in general, however, do not have emotions. As McDermott (2011, p. 2) states, “ethical behavior is an extremely difficult area to automate, both because it requires ‘solving all of AI’ and because even that might not be sufficient.” We are still far from having such capable machines.

AI ethics is not a new thing, nor did it start with the use of generative AI in the 2020s. Early debates on AI ethics, though not known under that name then, are as old as the AI field itself. Take, for instance, Norbert Wiener’s concerns. “As machines learn, they may develop unforeseen

strategies at rates that baffle their programmers. [...] We must always exert the full strength of our imagination to examine where the full use of our new modalities may lead us" (Wiener, 1960, p. 1355). "[W]hat, for example, are the scientist's responsibilities with respect to making his work public? And to whom (or what) is the scientist responsible?" (Weizenbaum, 1976, p. 8).

Ethical concerns regarding the use of technologies are not new in education either:

Many critics of education are impatient with the deficiency of our schools. They decry the tendency for education to lag far behind industry in automation. [...] But we must bear in mind that while the product of industry is an automobile, a refrigerator, or a washing machine, the product of education is a human being. (Watters, 2021, p. 258, citing Daniel Tanner, 1964)

We must consider the odds of the programs and systems that pupils, students, teachers, and administrators use or plan to use in educational contexts, for the harms and negative implications are not an imagined illusion but a real threat and a constant concern (see above and also Holmes et al., 2025; Huang, 2023; Wieczorek et al., 2025. See also below). This is why "ethical AI principles" are not "a specific political agenda," as Downes (2026, p. 2) mentions with apparent disdain, but *are* political in their entirety, as AI and technology *also are* (Bartoletti, 2020; Coeckelbergh, 2022; Crawford, 2022; McQuillan, 2022). "AI ethics is about technological change and its impact on individual lives, but also about transformations in society and in the economy" (Coeckelbergh, 2020, p. 7), which are politically loaded and determined. "AI ethics is about the lives of people and it is about policy" (Coeckelbergh, 2020, p. 62).

### 3.2 Which ethical principles? On methodology

In his Section 2.1, Downes (2026, p. 3) lists ethical principles found with Google Search's (generative) AI. This has serious methodological issues. For instance, the terms Downes used in his Google search are not specified, and none of the results he lists cite the original sources, which might be contributing to plagiarism or intellectual property issues (Appel et al., 2023; Crawford & Schultz, 2024; Panwar, 2025). As is well known, the output of generative AI can never be trusted: it is ripe with falsities, misinformation, and poor-quality content (Harding, 2024; Huang et al., 2024; Wiggers, 2024), and generative AI is not adequate for information seeking (Shah & Bender, 2022; 2024), among other reasons.

There is already a myriad of well-known works and policies that deal with ethical principles. Thus, any AI-generated output without critical evaluation of its origin or proper attribution must be avoided. Fundamentally, relying on AI-generated output obscures or erases the provenance of the contributions, does not credit the original authors, thereby making them and their work untraceable, undermines intellectual authorship and epistemic fairness, removes critical context, and misleads readers about the origin and credit of the sources (Earp et al., 2025). Furthermore, credit is a pillar of research integrity and is inherently intertwined with accountability and transparency (Kiermer et al., 2026; Nature, 2026).

Take, for example, how the AI HLEG (2019a, p. 14) defines ethical principles in the context of AI systems based on fundamental rights. Also, how those ethical principles are then "translated into concrete *requirements* to achieve *trustworthy AI* [and] are applicable to different stakeholders partaking in AI systems' life cycle: developers, deployers and end-users, as well as the broader society" (emphasis added), all of them relevant in educational settings and the majority overlapping with Downes's (2026) AI-generated list. Not giving credit to already existing works on those very topics risks reinventing them anew.

Jobin et al. (2019, p. 391)—in a widely cited paper with over 82k accesses at the time of writing this commentary—analysed "84 documents containing ethical principles or guidelines for AI".

They found that “[n]o single ethical principle appeared to be common to the entire corpus of documents, although there is an emerging convergence around the following principles: transparency, justice and fairness, non-maleficence, responsibility, and privacy” (p. 391). Attard-Frost et al. (2022) reviewed 47 ethics guidelines relevant to business practices and found, similarly, principles for AI ethics that are key to the ethics of AI business practices. Since AI tools and systems used in education are developed by AI companies, all those principles are relevant to them too, as the business practices behind them are. The work of those authors would be obscured if not properly cited.

State-of-the-art literature, frameworks, policies, analyses, and guidelines on ethical AI principles should have been considered to extract relevant ones from them. See also UNESCO (2022) and Wieczorek et al. (2025), as well as both the *Unified Framework of Five Principles for Ethical AI* (Floridi & Cows, 2019)—Downes briefly comments it—and the other sources from which these other authors extracted ethical principles from. Referring to AI’s output, i.e., to generative “AI” results, as the main source for an analysis about ethical principles erodes scientific integrity.

Moreover, in his Section 2.2, Downes (2026) only mentions some possible or potential benefits of the application of AI in education, part of which are not such but hype, by analysing only one source, i.e., Cardona et al. (2023). However, no single issue, problem, or harm, among the many that exist, is commented on. We refer the reader to Abbas et al. (2024), Bastani et al. (2024), Bauer et al. (2025), Fergusson et al. (2024), Forbes and Guest (2025), Guest et al. (2025), Laird et al. (2025), Lin et al. (2023), Weidlich et al. (2025), and Williamson et al. (2024), to name but a few works, to know about the dark side of “AI” in education, including empirical findings. These topics are a huge omission in Downes (2026).

### 3.3 The illusion of ethical use out of unethical practices and systems

It is worth mentioning that ethical principles and related requirements are not meant to be coded into or fulfilled by AI systems or algorithms independently, but are principles and requirements *humans* must respect and ensure (i.e., implement and evaluate) throughout the entire AI system’s life cycle. For example, programs cannot be held accountable; humans are. This is why the AI HLEG of the European Commission expects all interested groups to fulfil concrete roles depending on their types of interactions with AI systems.

[D]evelopers should implement and apply the requirements to design and development processes; [d]eployers should ensure that the systems they use and the products and services they offer meet the requirements; [and e]nd-users and the broader society should be informed about these requirements and able to request that they are upheld. (AI HLEG, 2019a, p. 11)

These good practices are, however, often ignored. On the other hand,

An [AI system] cannot genuinely be deemed ‘ethical’ without accounting for the business it is involved in. Perhaps the components involved in the system’s algorithmic decision-making—e.g., its data inputs, classification categories, and model characteristics—will operate ethically, but there is no guarantee that the system’s operation will remain ethical if its scope is widened to account for its business context. (Attard-Frost et al., 2022, p. 390)

*AI business practices* are defined as “the iterative political and economic behaviours involved in the organized resourcing, design, development, deployment, and use of an [AI system]” (Attard-Frost et al., 2022, p. 390), and therefore,

AI ethics research requires a much sharper focus on issues of business practice. [W]ithout more robust analysis of the business practices and political economies involved in [AI

systems'] development and use, AI ethics guidelines will remain dangerously limited in their ability to hold powerful businesses to account for their unethical practices. *There can be no ethical AI without ethical businesses to build it* (Attard-Frost et al., 2022, p. 397; emphasis added).

Artificial intelligence technologies, the ones Downes (2026) refers to in his paper—though it seems he only refers to one very special case of technique among the many that exist under the AI umbrella, especially the data and resources-intensive ones, are unethical at their core. They are based on unethical practices that cause harm, as it has been extensively documented in the literature, including factual and empirical research-based literature (Bender et al., 2021; Birhane et al., 2024; Dang & Liu, 2025; Eubanks, 2019; Fergusson et al., 2024; Goetze, 2024; Greenbaum & Gerstein, 2025; Hao & Seetharaman, 2023; Jargon & Schechner, 2025; Luccioni et al., 2024; Mejías & Coudry, 2024; Muldoon et al., 2024; Panwar, 2025; Vassel et al., 2024; Walther, 2024; Wilkins, 2025).

We posit that, thus, there can be no “ethical AI” coming out of the application of unethical business practices, and, as a consequence, there can be no possibility to make them ethical *a posteriori* or by just “acknowledging” the harms (Guersenzvaig & Monett, 2026; Guest et al., 2025; Monett & Paquet, 2025). That would be ethics washing and critical washing (Gerdes, 2022; Guest et al., 2025; Metzinger, 2019; Ochigame, 2019; van Maanen, 2022). In educational settings or contexts where such technologies are aggressively targeting users (teachers, students, administrators, and parents), for instance, they are predestined to harm. AI systems not only harms educational ecosystems but also students in particular, including their learning, cognitive abilities, long-term retention of content, career chances, and academic integrity, among other things (see Barcaui, 2025; Bastani et al., 2024; Fan et al., 2024; Gerlich, 2025; Guersenzvaig et al., 2025; Köbis et al., 2025; Laird et al., 2025; Lee et al., 2025; Liang et al., 2023; Lodge & Loble, 2026; Stadler et al., 2024; Watermeyer et al., 2023; Williamson et al., 2024; and Wong, 2025 for evidence about those harms). No current or potential benefit or use of a technology absolves it, or its creators, from the unethical practices that sustain it and the harms they account for.

Without an account for the unethical business practices underneath the AI applications that are being used (or planned to be used or deployed) in education, or for their limitations and the harms they inflict or might cause, any ethical qualifier applied to those AI products is destined to obfuscate their real nature, the purposes of their vendors, and their negative consequences.

## 4 A rebuttal to Downes’ analysis of ethical principles

This section will focus on those topics and/or arguments in Downes (2026) that are controversial. Comments will not be made on those that make sense.

### 4.1 On fairness (Downes, 2026, Sect. 3.1)

It is not true that fairness in AI only “amounts to the desire to, from a position of privilege, set those parameters that define where we will be ‘fair’ and where we will make actual ethical decisions” (Downes, 2026, p. 5). That is a very narrow definition of how fairness is considered in AI (Caton & Haas, 2024). Parameters are not the only elements of algorithms that are set when coding them, or that may play a role in fairness-related AI-based computational approaches. AI models can also deal with sensitive variables, for instance, and it can be that the treatment of sensitive variables is strongly related, e.g., to privacy issues that go beyond disadvantaged people and concern all kinds of targeted subjects in an analysis. What is more, some fairness-related approaches focus on the direct algorithmic impact on model accuracy and generalisability, rather than on actual ethical decisions that may derive from its use. On the other hand, data alone (or

the absence of it) might contribute to unfair decisions, for which parameter settings or variables use—i.e. from the perspective of the algorithm and its coder—would be insufficient. Also, because in some works “the notion of ‘fairness’ is not across multiple users or agents but across multiple objectives or criteria” (Reuel & Ma, 2024) that are optimised and don’t involve ethical decisions even indirectly.

#### 4.2 On transparency and explainability (Downes, 2026, Sect. 3.2)

Downes (2026) argues that “[t]he intuition behind [the transparency] condition is that AI systems should be designed in such a way that allows users to understand how they work” (p. 5). The goal is not to know how those systems work, because, according to Downes, laypeople might not understand the related technicalities, but how a decision was made, i.e., which were the conditions or parameters (e.g. whether the gender or the social group are conducing to unjust associations) that led to that decision, something especially relevant when contesting wrongdoing should the systems’ output discriminate or harm. In other words, it is not a matter of which single algorithmic operations or weight values the connections of a neural network have that needs to be known, but how the models deal with parameters and what they are, like, for example, if the place where the datafied subjects live or any information related to it that determines a wrong decision. The objective is not to comprehend or track billions of bits and the machine operations on them; posed this way, the problem is overcomplicated and exaggerated unnecessarily.

Furthermore, Downes’s (2026) comment on *conterfactuals* (p. 6) as a proxy for explainability is a well-known topic in the explainability community, also challenging technologically. Yet, the analysis of those other possible outcomes when dealing with conterfactuals is something machines cannot do unless programmed to do so, at least with current technologies. There exist several theories and methods for *deductive* and *inductive* reasoning, but none for *abductive* reasoning that satisfactorily explains or deals with impossibilities and common sense, however, a significant but missing piece in today’s AI-based decision making and steppingstone for genuine intelligence in machines (Cantwell Smith, 2019). Counterfactuals *are* part of the path in this direction (Pearl, 2018; Verma et al., 2024), yet an open research field with almost no or very few practical applications worldwide.

#### 4.3 On accountability (Downes, 2026, Sect. 3.3)

Downes (2026) argues that accountability “would help ensure that AI is used responsibly and that those responsible can be held liable for any negative consequences” (p. 6), which is only part of this ethical principle’s goal, though. What matters is not only to guarantee the correct use—what people would do when using or with AI, but most importantly to ensure that, in the case of harm, the ones who designed or developed the system or made wrongful decisions are held accountable. Downes (2026) only analyses the users and none of the other people involved in the AI life cycle.

He then goes on and compares the use of AI with using a hammer or a gun, saying that we are expecting its behaviour to be ethical by centring the discussion only on the technology and what can be done with it. But AI is not “just” a tool (Laba, 2025). And algorithms cannot be held accountable: it has always been about *people* and not about the technology. Downes’s (2026) analysis misses the point that what is called “AI” is not only a concrete software program, but a huge socio-technical, political, environmental, and cultural complex system of artefacts, relationships, dynamics, and decisions made all the way down by a few unelected people to cement profit, control, and power, and which affects entire generations, societies, ecosystems, and the planet (Crawford, 2022; Guest et al., 2025; Hao, 2025; McQuillan, 2022). This is why it is,

*of course*, much needed and utterly important to know *who* is accountable, contrary to what Downes suggests.

#### 4.4 On privacy and data protection (Downes, 2026, Sect. 3.4)

The AI HLEG alerts about the tensions that can exist between ethical principles:

[I]n various application domains, *the principle of prevention of harm* and *the principle of human autonomy* may be in conflict. Consider as an example the use of AI systems for ‘predictive policing’, which may help to reduce crime, but in ways that entail surveillance activities that impinge on individual liberty and privacy. [...] There may be situations, however, where no ethically acceptable trade-offs can be identified. Certain fundamental rights and correlated principles are absolute and cannot be subject to a balancing exercise (e.g. human dignity). (AI HLEG, 2019, p. 13; emphasis in the original)

Yet, Downes (2026) generalises and poses an unsolvable dichotomy, claiming that “[w]e can have accountability, or we can have data protection, but we can't have both at the same time for the same data” (p. 8). Whilst this can be true in some situations, there are others where we can have *both*. On the one hand, accountability is not only about who uses a system, as explained above; on the other, privacy and data protection can be safeguarded, for example, in a social media platform, but at the same time with strict measures in place that make people accountable for prohibited behaviours without necessarily revealing their data to third parties.

Later on, Downes (2026) succumbs to the AI hype, again, saying that “[a]n AI, for example, can recognize a person by their gait. Humans can do this too, but only for people they already know well; an AI can do it for anyone” (p. 8). No, “an” AI cannot do it for anyone; this statement is misleading and ignores how the technology works. The contrary has been demonstrated more than once: AI systems are deployed and put in the hands of police departments, to name an example of an application domain, that misclassify and discriminate against people by matching them to collected data from other unrelated people the algorithms were trained on. Mathematical operations on pixels and image sequences still cannot generalize to distinguish *any* people based on limited traits. As expected, the most affected and abused individuals and groups are always the already marginalised and disadvantaged (Kalluri et al., 2025).

#### 4.5 On inclusiveness (Downes, 2026, Sect. 3.5)

How Downes (2026) refers to the ethical principle of inclusiveness is also problematic. His position is that *if this is not inclusive here, then it makes no sense to have inclusiveness there; thus, some people don't deserve it*. For example, he notes, “Almost nothing we produce is required to be used to benefit all members of society. It's hard to see how inclusiveness stands as an ethical value specifically for AI” (Downes, 2026, p. 9). Again, it is not that AI is something special that deserves attention now and not before, but that AI-based socio-technical systems and the decisions that result from them amplify, propagate, extend, and perpetuate the biases, exclusiveness, and harm not only present in the data, but also coming from those who design, implement, test, deploy, control, and use those systems (Birhane, 2021; Birhane et al., 2024; Eubanks, 2019; Fergusson et al., 2024).

Other parts of the section on inclusiveness are as concerning as defeatist and deterministic. For instance, Downes (2026) maintains that “the principle of diversity, equity, and inclusion (DEI) has been explicitly rejected by the United States government” (p. 9) and adds nothing else to the discussion. Not questioning its discriminative and abusive nature in a section dedicated to the opposite goal sides with those in power, who are not interested in the big majority's well-being. *It is so; there is nothing to do*, seems to be the message. Excluding people and optimizing for certain

traits has been ingrained in and related to the AI field for many years (Geburu & Torres, 2024), an unacceptable condition to leave uncriticised.

#### 4.6 On reliability and safety (Downes, 2026, Sect. 3.6)

The next section goes similarly, with quite the same position as in previous ones: *this is already bad here; thus, there is no sense in not wanting badness there*, now referred to as reliability and safety. According to Downes (2026), there is no need to do anything special in the case of AI in matters of reliability and safety, the technology that amplifies the harms produced by unreliable and insecure systems, because that's not something people value as having more priority in their lives.

Reliability and safety are only partially analysed. A reliable system is one where one can expect the same answer the next time it is used with the same input, for instance. One can trust the output because it doesn't change, or because the sources the output is based on are reliable and cited reliably. Similarly, a safe system is safe, and thus robust, against attacks that leak the data it works with; it doesn't leave backdoors for third parties to exploit. None of this is analysed in Downes (2026), who focuses only on the perspective of the final user. For the sake of argument, or to exemplify, Downes (2026) refers only to technical aspects when analysing reliability and safety. Expressions like "Computers are reliable [...]. Our computers are safe" (Downes, 2026, p. 10) again generalize characteristics of digital systems that are not true in the normal case. Computers, digital systems, algorithms, and AI are neither reliable nor safe: they are hacked and the hardware fails often, mainly because of bad design decisions. The technology cannot be deemed to be reliable or safe: for example, injecting damaging prompts into chatbots' conversations is as frequent as the software band aids to avoid them (Claburn, 2026). It is impossible to prevent those systems from unreliability or vulnerability completely. This doesn't mean we should accept unethical practices regarding reliability and safety!

Downes (2026) goes on to make the assumption that "society as a rule tends to balance considerable levels of risk against anticipated benefits and trust in the provider" (p. 10). Actually, and especially in the case of AI, what humans are good at is gullibility, myths, illusioned machine capabilities, and wishful thinking, blatantly failing to consider the harms and risks of current systems, but blindly believing in what vendors, boosters, and techno-oligarchs sell as "good AI" or "for the good of humanity" (Adib-Moghaddam, 2025; Greenspan, 2009; Madianou, 2024; McDermott, 1976; Weizenbaum, 1976).

Downes (2026, p. 10) then argues that "[t]he suggestion that we should change this approach now, for this technology, does not appear to be any longstanding ethical principle but rather a latent conservatism that mistrusts new technologies in general, in other words, neo-Luddites (Lamont, 2024)." This doesn't deserve time falsifying, but demonstrates ignorance about what Luddism actually was and is: it is about mistrusting, refusing, and opposing those who use technologies to exploit, extract from, control, and abuse other humans, societies at large, ecosystems, and the planet, to exercise domination, conserve unaccountability, and increase profit (Binfield, 2015; Merchant, 2023; Mueller, 2021; Sadowski, 2025). It is not about AI, it is about who and what they do to us with "AI."

#### 4.7 On human-centred design (Downes, 2026, Sect. 3.7)

The section on human-centred design is not better. Downes (2026) cites problematic quotes (p. 11) using anthropomorphic (e.g., "We also need to design and develop systems with [...] possibly even the needs of artificial intelligences in mind"), techno-determinist (e.g., "Whatever we build, we will adapt to it"), and techno-solutionist language (e.g. "the imperatives and the affordances of human and machine shape each other"), as also absurd examples (e.g., "The needs of humans are over-ruled by machines all the time. The ATM won't let me spend more money than I have"),

to make a point by underestimating the importance of human-centred design. What Downes (2026) fails to acknowledge is that AI is neither artificial nor intelligent: it is very human at its core; it cannot exist nor function, especially data-intensive AI, without the digital theft of both human intellect and works, mostly without consent, as well as labour in the order of millions worldwide, e.g. for gathering the data, cleaning and pre-processing it, annotating it, moderating it, giving feedback about it, and so on (Crawford, 2022; Hao & Seetharaman, 2023; Koebler, 2026; Mejías & Couldry, 2024; Muldoon et al., 2024). AI is demonstrably very stupid too (Bishop, 2021) and far from being “intelligent” as humans are (see Section 2 above).

As Guest (2025, p. 3) reminds us,

AI is human-centric, *not* because it behaves like or is designed to be like humans, but because it requires a ghost in the machine, often literally an obfuscated human-in-the-loop to properly function (also see Guest & Martin, 2025) because AI *is* humans albeit in fetishised (Braune, 2020; Morris, 2017; Mota & Cosentino Filho, 2024), obfuscated forms (e.g., Erscoi et al., 2023).

Humans—the ones who are involved in the development of, who use, and whose lives and future are being affected by AI-based systems—and their needs should always be at the centre and in all phases of any software and the socio-technical systems it is part of. In some cases, even not only humans but broader ecosystems. That’s what this ethical principle is for.

#### 4.8 On non-maleficence (Downes, 2026, Sect. 3.8)

The general tone of this section is similar to the former ones: if technology does exist that causes harm, and because some are clear in its nature and we avoid them, or because there are many harms already, including “necessary” and “allowable” harms, then the principle of non-maleficence is kind of superfluous, especially if the law prevents its application involving stakeholders affected by it. Not humans as targets of systems of oppressions and harms, but stakeholders, the ones profiting from the harms the systems produce. And we had better do nothing against something inevitable and determined, don’t we!

The contrary *is* precisely what this ethical principle and regulation (like the European AI Act, among others) are for: neither stakeholders’ profit nor technologists’ excesses can ever be more important than the humans who are harmed by the technology they gamble with. Human rights and human lives are above any AI-based artefact and the monetary purposes of its enablers.

#### 4.9 On sustainability (Downes, 2026, Sect. 3.9)

The last section deals with the principle of sustainability. It might be that Downes is not aware of the literature covering the tremendous negative impact of AI-based technologies for the environment, communities, and the planet; otherwise, he may not have made the statement that “[t]he power consumed by AI is not even a rounding error when compared to the total amount of energy consumed by humans. It’s far too small to be considered even that” (Downes, 2026, p. 12). It seems his observation needs urgent revision. Take, for instance:

The carbon footprint of AI systems alone could be between 32.6 and 79.7 million tons of CO<sub>2</sub> emissions in 2025, while the water footprint could reach 312.5–764.6 billion L. To put this into perspective, this is in the same range as the carbon footprint of New York City (52.2 million tons of CO<sub>2</sub> emissions in 2023). Similarly, the water footprint of AI systems may be in the same range as the entire global annual consumption of bottled water (446 billion L). (de Vries-Gao, 2026, p. 9)

And a report of the International Energy Agency estimates that “data centres’ total electricity consumption could reach more than 1000 TWh in 2026 [...] roughly equivalent to the electricity

consumption of Japan” (Çam et al., 2024). Other studies report similar concerns. (Crawford, 2024; de Vries, 2023; Garg & Kitsara, 2025; Niranjana, 2026; O'Donnell & Crownhart, 2025; Xiao et al., 2025, Zewe, 2025), as also the extraction of basic natural resources and occupation of land to build new data centres (Jiménez Arandia et al., 2025; Mejías & Couldry, 2024).

Even more concerning is how Downes (2026) compares “the energy we’re consuming if we’re not using AI” with “how much energy and resources it takes to pay for and support 600,000 translators” (p. 12), thereby suggesting that the work of those 600,000 translators, with which they sustain themselves and their families, is not worth their value and work and, thus, should better be replaced by AI. Yet, even after around 80 years of automatic translation’s modern history, AI is well-known to remain far behind regarding the nuanced, highly specialized, context-sensitive, culturally-dependent, and meaning-accurate human touch, something that machines are far to equate, needless to say, surpass (Naveen & Trojovský, 2024; Pang et al., 2025). A similar, recent comment by Sam Altman, the CEO of OpenAI, questions, in an absurd comparison, the energy consumption of raising a child and the one needed to train AI (Berger, 2026). Such excuses obfuscate the real nature, resources, infrastructure, algorithm, and data processes, data itself, as well as the *millions* of humans historically building and feeding “AI” for it to work. Such category errors and the issues they hide are inexcusable.

## 5 Conclusion

Downes’s (2026) use of Google Search’s (generative) AI felt short of relevance. It left out works and ethical principles essential to any digital—in general, and AI, in particular—technology of our times. Ethical principles like human agency, accuracy, reproducibility, societal well-being, non-discrimination, and data governance, among others, were brushed away from the AI ethics spectrum with an “AI” click. If some of the ethical principles analysed in Downes’s paper were considered to be unnecessary or superfluous, what could be expected for the ones that were left out?

AI ethics, digital ethics, ethical principles, and so on are far from being perfect abstractions to guide human decision-making and technology developments, but in times where ethics washing and critical washing (Gerdes 2022; Guest et al., 2025; Metzinger, 2019; Ochigame 2019; van Maanen, 2022) tend to dominate politics and governance, any sound criticism is more than welcome. However, the discourse and validity of ethical principles related to AI should start from a solid analysis that acknowledges their importance, for example, in

ensuring a high level of protection of health, safety, fundamental rights as enshrined in the Charter of Fundamental Rights of the European Union (the ‘Charter’), including democracy, the rule of law and environmental protection, to protect against the harmful effects of AI systems in the Union, and to support innovation. (EU, 2024)

Contrary to what Downes (2026) thinks, that is, university professors and the institutions they represent “are frequently among those who define and recommend adherence to ethical principles in all we do, including the deployment of AI” (p. 13), there are many members in ethical commissions, companies, institutions, and organisations worldwide who are not academics. Take a look, for example, at the 52 members of the AI HLEG<sup>4</sup> of the European Commission, experts who developed foundational materials key to the first AI regulation of its type worldwide, the AI Act of the European Union.

---

<sup>4</sup> See <https://ec.europa.eu/futurium/en/european-ai-alliance/ai-hleg-steering-group-european-ai-alliance.html>

Instead of downplaying the role of academia, as Downes (2026) does when he claims academics' efforts to define and recommend adherence to ethical principles are "laudable, but it's hard not to observe a certain amount of misdirected energy on the part of academia" (p. 13), we would like to call for a radically different approach:

The role of an intellectual should consist not so much in instructing the masses, to the point of preaching to them from their ivory tower, as in understanding the phenomena that are developing, in order to sound the alarm with arguments and awareness when necessary. (Sadin, 2023, p. 106; translated from Spanish)

Any amount of energy invested in the latter is more than welcome. It is much more: it is the *raison d'être* as educators, academics, and citizens; it is vital.

## References

- Abbas, M., Jam, F. A., & Khan, T. I. (2024). Is it harmful or helpful? Examining the causes and consequences of generative AI usage among university students. *International Journal of Educational Technology in Higher Education*, 21, 10. <https://doi.org/10.1186/s41239-024-00444-7>
- Adib-Moghaddam, A. (2025). *The Myth of Good AI: A manifesto for critical Artificial Intelligence*. Manchester University Press.
- AI HLEG (2019a). *Ethics Guidelines for Trustworthy AI*. High-Level Expert Group on Artificial Intelligence, European Commission. [https://ec.europa.eu/newsroom/dae/document.cfm?doc\\_id=60419](https://ec.europa.eu/newsroom/dae/document.cfm?doc_id=60419)
- AI HLEG (2019b). *A Definition of AI: Main Capabilities and Disciplines*. High-Level Expert Group on Artificial Intelligence, European Commission. [https://ec.europa.eu/newsroom/dae/document.cfm?doc\\_id=60651](https://ec.europa.eu/newsroom/dae/document.cfm?doc_id=60651)
- Ali, S.M., Dick, S., Dillon, S., Jones, M.L., Penn, J., & Staley, R. (2023). Histories of artificial intelligence: A genealogy of power. *BJHS Themes*, 8, 1–18. <https://doi.org/10.1017/bjt.2023.15>
- Appel, G., Neelbauer, J., & Schweidel, D. A. (2023, April 7). Generative AI has an intellectual property problem. *Harvard Business Review*. <https://hbr.org/2023/04/generative-ai-has-an-intellectual-property-problem>
- Attard-Frost, B., De los Ríos, A., & Walters, D. R. (2022). The ethics of AI business practices: A review of 47 AI ethics guidelines. *AI and Ethics*, 3, 389–406. <https://doi.org/10.1007/s43681-022-00156-6>
- Barcaui, A. (2025). ChatGPT as a cognitive crutch: Evidence from a randomized controlled trial on knowledge retention. *Social Sciences & Humanities Open*, 12, 102287. <https://doi.org/10.1016/j.ssaho.2025.102287>
- Barrow, N. (2024). Anthropomorphism and AI hype. *AI Ethics*, 4, 707–711. <https://doi.org/10.1007/s43681-024-00454-1>
- Bartneck, C., Lütge, C., Wagner, A., Welsh, S. (2019). *Ethics in AI and Robotics* (Original: Ethik in KI und Robotik). Carl Hanser Verlag.
- Bartoletti, I. (2020). *An artificial revolution: On power, politics and AI*. The Indigo Press.
- Bastani, H., Bastani, O., Sungu, A., Ge, H., Kabakci, Ö, & Mariman, R. (2024). Generative AI can harm learning. *The Wharton School Research Paper*, 1–59. <http://dx.doi.org/10.2139/ssrn.4895486>
- Bauer, E., Greiff, S., Graesser, A. C., Scheiter, K., & Sailer, M. (2025). Looking beyond the hype: Understanding the effects of AI on learning. *Educational Psychology Review*, 37, 45. <https://doi.org/10.1007/s10648-025-10020-8>
- Bender, E.M. (2024, April 2). Doing their hype for them: Defeatist, second-hand hype goes to college. *Mystery AI Hype Theater 3000: The Newsletter*. <https://buttondown.email/maiht3k/archive/doing-their-hype-for-them/>

- Bender, E. M., Gebru, T., McMillan-Major, A., & Shmitchell, S. (2021). On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '21, 610–623. <https://doi.org/10.1145/3442188.344592>
- Bender, E.M., & Hanna, A. (2025). *The AI Con: How to Fight Big Tech's Hype and Create the Future We Want*. Harper.
- Berger, E. (2026). Sam Altman defends AI's energy toll by saying it also takes a lot to 'train a human.' *The Guardian*. <https://www.theguardian.com/technology/2026/feb/23/sam-altman-openai-energy-use-datacenters>
- Binfield, K. (2015). *Writings of the Luddites*. Johns Hopkins University Press.
- Birhane, A. (2021). Algorithmic injustice: a relational ethics approach. *Patterns*, 2(2), 100205. <https://doi.org/10.1016/j.patter.2021.100205>
- Birhane, A., Dehdashtian, S., Prabhu, V., Boddeti, V. (2024). The Dark Side of Dataset Scaling: Evaluating Racial Classification in Multimodal Models. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '24, 1229–1244. Association for Computing Machinery. <https://doi.org/10.1145/3630106.3658968>
- Bishop, J. M. (2021). Artificial intelligence is stupid and causal reasoning will not fix it. *Frontiers in Psychology*, 11, 513474. <https://doi.org/10.3389/fpsyg.2020.513474>
- Çam, E., Hungerford, Z., Schoch, N., Pinto Miranda, F., & Yáñez de León, C.D. (2024). Electricity 2024: Analysis and forecast to 2026. *International Energy Agency*. <https://www.iea.org/reports/electricity-2024>
- Cantwell Smith, B. (2019). *The promise of artificial intelligence: Reckoning and judgement*. The MIT Press.
- Caton, S., & Haas, C. (2024). Fairness in machine learning: A survey. *ACM Computing Surveys*, 56(7), 166, 1–38. <https://doi.org/10.1145/3616865>
- Cicurel, R., & Nicolescu, M.A.L. (2015). *The relativistic brain: How it works and why it cannot be simulated by a Turing machine*. Kios Press.
- Claburn, T. (2026). OpenAI putting bandaids on bandaids as prompt injection problems keep festering. *The Register*. [https://www.theregister.com/2026/01/08/openai\\_chatgpt\\_prompt\\_injection/](https://www.theregister.com/2026/01/08/openai_chatgpt_prompt_injection/)
- Coeckelbergh, M. (2020). *AI ethics*. The MIT Press.
- Coeckelbergh, M. (2022). *The political philosophy of AI: An introduction*. Polity.
- Crawford, K. (2022). *Atlas of AI: Power, politics, and the planetary costs of artificial intelligence*. Yale University Press.
- Crawford, K. (2024, February 20). Generative AI's environmental costs are soaring — and mostly secret. *Nature*. <https://www.nature.com/articles/d41586-024-00478-x>
- Crawford, K., & Schultz, J. (2024, January 16). Generative AI is a crisis for copyright law. *Issues in Science and Technology*, 79–80. <https://issues.org/generative-ai-copyright-law-crawford-schultz/>
- Critical AI (n.d.). The AI hype wall of shame. *Critical AI*. <https://criticalai.org/the-ai-hype-wall-of-shame/>
- Dang, J., & Liu, L. (2025). Dehumanization risks associated with artificial intelligence use. *American Psychologist*. <https://doi.org/10.1037/amp0001542>
- De Vries, A. (2023). The growing energy footprint of artificial intelligence. *Joule*, 7(10), 2191–2194. <https://doi.org/10.1016/j.joule.2023.09.004>
- De Vries-Gao, A. (2026). The carbon and water footprints of data centers and what this could mean for artificial intelligence. *Patterns*, 7(1), 101430. <https://doi.org/10.1016/j.patter.2025.101430>
- Doroudi, S. (2023). The intertwined histories of artificial intelligence and education. *International Journal of Artificial Intelligence in Education*, 33, 885–928. <https://doi.org/10.1007/s40593-022-00313-2>

- Downes, S. (2026). On ethical AI principles. *Journal of Open, Distance, and Digital Education*, 3(1), 1–18. <https://doi.org/10.25619/622a7242>
- Duarte, T., Barrow, N., Bakayeva, M., & Smith, P. (2024). Editorial: The ethical implications of AI hype. *AI Ethics*, 4, 649–651. <https://doi.org/10.1007/s43681-024-00539-x>
- Earp, B.D., Yuan, H., Koplin, J., Porsdam Mann, S. (2025). LLM use in scholarly writing poses a provenance problem. *Nature Machine Intelligence*, 7, 1889–1890. <https://doi.org/10.1038/s42256-025-01159-8>
- EU (2024). Artificial intelligence act. Regulation (EU) 2024/1689 of the European Parliament and of the Council of 13 June 2024. *European Parliament & Council of the European Union*. <http://data.europa.eu/eli/reg/2024/1689/oj>
- Eubanks, V. (2019). *Automating inequality: How high-tech tools profile, police, and punish the poor*. St. Martin's Press.
- Fan, Y., Tang, L., Le, H., Shen, K., Tan, S., et al. (2024). Beware of metacognitive laziness: Effects of generative artificial intelligence on learning motivation, processes, and performance. *British Journal of Educational Technology*, 56(2), 489–530. <https://doi.org/10.1111/bjet.13544>
- Fergusson, G., Geoghegan, S., Schroeder, C., & Villegas Bravo, M. (2024). Generating harms: Generative AI's new & continued impacts. *Electronic Privacy Information Center*. <https://epic.org/wp-content/uploads/2024/05/EPIC-Generative-AI-II-Report-May2024-1.pdf>
- Floridi, L., & Cows, J. (2019, July 2). A unified framework of five principles for AI in society. *Harvard Data Science Review*, 1(1). <https://doi.org/10.1162/99608f92.8cd550d1>
- Floridi, L., & Nobre, A. C. (2024). Anthropomorphising machines and computerising minds: The crosswiring of languages between artificial intelligence and brain & cognitive sciences. *Minds and Machines*, 34(5). <https://doi.org/10.1007/s11023-024-09670-4>
- Forbes, S.H., & Guest, O. (2025). To improve literacy, improve equality in education, not large language models. *Cognitive Science*, 49(4), e70058. <https://doi.org/10.1111/cogs.70058>
- Garg, A., & Kitsara, I. (2025, February 26). The hidden cost of AI: Unpacking its energy and water footprint. *OECD. AI Policy Observatory*. <https://oecd.ai/en/wonk/the-hidden-cost-of-ai-energy-and-water-footprint>
- Geburu, T., & Torres, É.P. (2024). The TESCREAL bundle: Eugenics and the promise of utopia through artificial general intelligence. *First Monday*, 29(4). <https://doi.org/10.5210/fm.v29i4.13636>
- Gerdes, A. (2022). The tech industry hijacking of the AI ethics research agenda and why we should reclaim it. *Discover Artificial Intelligence*, 2, 25. <https://doi.org/10.1007/s44163-022-00043-3>
- Gerlich, M. (2025). AI tools in society: Impacts on cognitive offloading and the future of critical thinking. *Societies*, 15(1), 6. <https://doi.org/10.3390/soc15010006>
- Goetze, T.S. (2024). AI art is theft: Labour, extraction, and exploitation, Or, on the dangers of stochastic pollocks. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency*, FAccT'24, 186–196. <https://doi.org/10.1145/3630106.3658898>
- Greenbaum, D., & Gerstein, M. (2025). Hidden human costs of AI. *Science*, 387, 32–32. <https://www.science.org/doi/10.1126/science.adu1541>
- Greenspan, S. (2009). *Annals of gullibility: Why we get duped and how to avoid it*. Praeger Publishers/Greenwood Publishing Group.
- Guersenzvaig, A., & Monett, D. (2026). Resisting enchantment and determinism: How to critically engage with AI university guidelines. *Zenodo*. <https://doi.org/10.5281/zenodo.18282338>
- Guersenzvaig, A., Sánchez-Monedero, J., Gopegui, B., & Picassó i Piquer, M. (2025). Critical thinking, teaching, and generative artificial intelligence. (Original: Inteligencia artificial generativa en la educación universitaria: la urgencia de una perspectiva crítica). *Daimon Revista Internacional de Filosofía*, 94. <http://dx.doi.org/10.6018/daimon.636721>
- Guest, O. (2025). What does 'Human-Centred AI' mean? arXiv:2507.19960v2 [cs.AI].

- <https://doi.org/10.48550/arXiv.2507.19960>
- Guest, O., Suárez, M., Müller, B., van Meerkerk, E., Oude Groote Beverborg, A., de Haan, R., Reyes Elizondo, A., Blokpoel, M., Scharfenberg, N., Kleinherenbrink, A., Camerino, I., Woensdregt, M., Monett, D., Brown, J., Avraamidou, L., Alenda-Demoutiez, J., Hermans, F., van Rooij, I. (2025). Against the uncritical adoption of 'AI' technologies in academia. *Zenodo*. <https://doi.org/10.5281/zenodo.17065099>
- Hao, K. (2025). *Empire of AI: Dreams and nightmares in Sam Altman's OpenAI*. Penguin Random House.
- Hao, K., & Seetharaman, D. (2023, July 24). Cleaning up ChatGPT takes heavy toll on human workers. *The Wall Street Journal*. <https://www.wsj.com/articles/chatgpt-openai-content-abusive-sexually-explicit-harassment-kenya-workers-on-human-workers-cf191483>
- Harding, X. (2024, July 10). When search engine's AI overviews are bad, they're really bad – What's the fix? *Mozilla Foundation*. <https://www.mozillafoundation.org/en/blog/ai-overview-google-search/>
- Holmes, W., Mouta, A., Hillman, V., Schiff, D., Laak, K., & et al. (2025). Critical studies of artificial intelligence and education: Putting a stake in the ground. *Sociology Education eJournal*, SSRN. <http://dx.doi.org/10.2139/ssrn.5391793>
- Huang, L. (2023). Ethics of artificial intelligence in education: Student privacy and data protection. *Science Insights Education Frontiers*, 16(2), 2577–2587. <https://doi.org/10.15354/sief.23.re202>
- Huang, L., Yu, W., Ma, W., Zhong, W., Feng, Z., Wang, H., Chen, Q., Peng, W., Feng, X., Qin, B., & Liu, T. (2024). A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *ACM Transactions on Information Systems*, 43(2), 1–45, 42. <https://doi.org/10.1145/3703155>
- Jargon, J., & Schechner, S. (2025, November 6). Seven lawsuits allege OpenAI encouraged suicide and harmful delusions. *The Wall Street Journal*. <https://www.wsj.com/tech/ai/seven-lawsuits-allege-openai-encouraged-suicide-and-harmful-delusions-25def1a3>
- Jiménez Arandia, P., Dib, D., Alarcón, M. (2025, August 8). The backyard of AI: A map of the 21st century gold rush. *El País*. <https://pulitzercenter.org/stories/backyard-ai-map-21st-century-gold-rush>
- Jobin, A., Ienca, M., & Vayena, E. (2019). The global landscape of AI ethics guidelines. *Nature Machine Intelligence*, 1(9), 389–399. <https://doi.org/10.1038/s42256-019-0088-2>
- Kalluri, P.R., Agnew, W., Cheng, M., Owens, K., Soldaini, L., & Birhane, A. (2025). Computer-vision research powers surveillance technology. *Nature*, 643, 73–79. <https://doi.org/10.1038/s41586-025-08972-6>
- Kambhampati, S. (2024, March 6). Can large language models reason and plan? *Annals of the New York Academy of Sciences*. <https://doi.org/10.1111/nyas.15125>
- Kambhampati, S., Stechly, K., & Valmeekam, K. (2025, April 12). (How) Do reasoning models reason? *Annals of the New York Academy of Sciences*. <https://doi.org/10.1111/nyas.15339>
- Kiermer, V. Adams, S., Bibbins-Domingo, K., Flores Bueso, Y., Jamieson, K.H., Heber, J., Hosseini, M., Marušić, A., Nielsen, B., Skipper, M., Swamy, G.K., Wolf, S.M. (2026). Creating a responsible authorship culture in science: Anchoring authorship practices in principles of transparency, credit, and accountability. *Proceedings of the National Academy of Sciences U.S.A.*, 123(12), e2531268123. <https://doi.org/10.1073/pnas.2531268123>
- Köbis, N., Rahwan, Z., Rilla, R., Supriyatno, B.I., Bersch, C., Ajaj, T., Bonnefon, J.-F., & Rahwan, I. (2025). Delegation to artificial intelligence can increase dishonest behaviour. *Nature*, 646, 126–134. <https://doi.org/10.1038/s41586-025-09505-x>
- Koebler, J. (2026, March 12). 'AI Is African Intelligence': The Workers Who Train AI Are Fighting Back. *404 Media*. <https://www.404media.co/ai-is-african-intelligence-the-workers-who-train-ai-are-fighting-back/>

- Laba, N. (2025). AI is not a tool. *AI & Society*. <https://doi.org/10.1007/s00146-025-02784-y>
- LaGrandeur, K. (2024). The consequences of AI hype. *AI Ethics*, 4, 653–656. <https://doi.org/10.1007/s43681-023-00352-y>
- Laird, E., Dwyer, M., & Quay-de-la-vallee, H. (2025). Hand in hand: Schools' embrace of AI connected to increased risks to students. *Center for Democracy and Technology*. <https://cdt.org/wp-content/uploads/2025/10/FINAL-CDT-2025-Hand-in-Hand-Polling-100225-accessible.pdf>
- Landgrebe, J., & Smith, B. (2022). *Why machines will never rule the world: Artificial intelligence without fear*. Routledge.
- Lee, H.-P., Sarkar, A., Tankelevitch, L., Drosos, I., Rintel, S., Banks, R., & Wilson, N. (2025). The impact of generative AI on critical thinking: Self-reported reductions in cognitive effort and confidence effects from a survey of knowledge workers. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*, 1121. <https://doi.org/10.1145/3706598.3713778>
- Lehman, J., Clune, J., & Risi, S. (2014). An anarchy of methods: Current trends in how intelligence is abstracted in AI. *IEEE Intelligent Systems*, 29, 56–52. <https://doi.ieeecomputersociety.org/10.1109/MIS.2014.92>
- Liang, W., Yuksekogonul, M., Mao, Y., Wu, E., Zou, J. (2023). GPT detectors are biased against non-native English writers. *Patterns*, 4(7), 100779. <https://doi.org/10.1016/j.patter.2023.100779>
- Lin, S. M., Chung, H. H., Chung, F. L., & Lan, Y. J. (2023). Concerns about using ChatGPT in education. In Huang, Y. M., Rocha, T. (Eds.), *Innovative Technologies and Learning*. Lecture Notes in Computer Science, 14099, 37–49. Springer, Cham. [https://doi.org/10.1007/978-3-031-40113-8\\_4](https://doi.org/10.1007/978-3-031-40113-8_4)
- Lipton, Z.C., & Steinhardt, J. (2019). Troubling trends in machine learning scholarship: Some ML papers suffer from flaws that could mislead the public and stymie future research. *Queue*, 17(1), 45–77. <https://doi.org/10.1145/3316774>
- Lodge, J.M., Loble, L. (2026). Artificial intelligence, cognitive offloading and implications for education. Report, 1–46, University of Technology Sydney, Australia. <https://doi.org/10.71741/4pyxmbnjq.31302475>
- Luccioni, S., Trevelin, B., & Mitchell, M. (2024, September 3). The environmental impacts of AI – Primer. *Hugging Face*. <https://huggingface.co/blog/sasha/ai-environment-primer#how-does-ai-use-energy>
- Madianou, M. (2024). *Technocolonialism: When technology for good is harmful*. Polity.
- Marchetti, A., Manzi, F., Riva, G., Gaggioli, A., & Massaro, D. (2025). Artificial intelligence and the illusion of understanding: A systematic review of theory of mind and large language models. *Cyberpsychology, Behavior, and Social Networking*, 28(7), 505–514. <https://doi.org/10.1089/cyber.2024.0536>
- Markelius, A., Wright, C., Kuiper, J., Delille, N., & Kuo, Y.-T. (2024). The mechanisms of AI hype and its planetary and social costs. *AI Ethics*, 4, 727–742. <https://doi.org/10.1007/s43681-024-00461-2>
- McCarthy, J., Minsky, M.L., Rochester, N., & Shannon, C.E. (1955). *A proposal for the Dartmouth summer research project on artificial intelligence*. <http://jmc.stanford.edu/articles/dartmouth/dartmouth.pdf>
- McDermott, D. (1976). Artificial intelligence meets natural stupidity. *ACM SIGART Bulletin*, 57, 4–9. <https://dl.acm.org/doi/10.1145/1045339.104534>
- McDermott, D. (2011). What matters to a machine? In S. Anderson, M. Anderson (Eds.), *Machine Ethics* (pp. 88–114). Cambridge University Press. <https://www.cs.yale.edu/homes/dvm/papers/whatmatters.pdf>
- McQuillan, D. (2022). *Resisting AI: An anti-fascist approach to artificial intelligence*. Bristol University Press.

- Mejías, U. A., & Couldry, N. (2024). *Data grab: The new colonialism of big tech and how to fight back*. The University of Chicago Press.
- Merchant, B. (2023). *Blood in the machine: The origins of the rebellion against big tech*. Little, Brown and Company.
- Metzinger, T. (2019, April 8). EU guidelines: Ethics washing made in Europe. *Tagesspiegel*. <https://www.tagesspiegel.de/politik/ethics-washing-made-in-europe-5937028.html>
- Monett, D., & Lewis, C. W. P. (2018). Getting clarity by defining artificial intelligence—A survey. In Müller, V. (eds), *Philosophy and Theory of Artificial Intelligence 2017*. PT-AI 2017. *Studies in Applied Philosophy, Epistemology and Rational Ethics*, 44, 212–214. Springer, Cham. [https://doi.org/10.1007/978-3-319-96448-5\\_21](https://doi.org/10.1007/978-3-319-96448-5_21)
- Monett, D., & Paquet, G. (2025). Against the commodification of education — if harms then not AI. *Journal of Open, Distance, and Digital Education*, 2(1), 1–24. <https://doi.org/10.25619/wazgw457>
- Mueller, G. (2021). *Breaking things at work: The luddites are right about why you hate your job*. Verso.
- Muldoon, J., Graham, M., & Cant, C. (2024). *Feeding the machine: The hidden human labor powering A.I.* Bloomsbury Publishing.
- Nature (2026). Credit in research goes hand in hand with responsibility. Editorial, *Nature*, 649, 527. <https://doi.org/10.1038/d41586-026-00006-z>
- Naveen, P., & Trojovský, P. (2024). Overview and challenges of machine translation for contextually appropriate translations. *iScience*, 27(10), 110878. <https://www.sciencedirect.com/science/article/pii/S2589004224021035>
- Nilsson, N.J. (2009). *The quest for artificial intelligence: A history of ideas and achievements*. Cambridge University Press. <http://www.cambridge.org/us/0521122937>
- Nilsson, N.J. (2012). John McCarthy 1927–2011: A biographical memoir. *Biographical Memoirs*, National Academy of Sciences, 1–17. <http://biographicalmemoirs.org/pdfs/mccarthy-john.pdf>
- Niranjan, A. (2026, January 3). ‘Just an unbelievable amount of pollution’: how big a threat is AI to the climate? *The Guardian*. <https://www.theguardian.com/technology/2026/jan/03/just-an-unbelievable-amount-of-pollution-how-big-a-threat-is-ai-to-the-climate>
- Ochigame, R. (2019, December 20). The Invention of ‘Ethical AI’: How Big Tech Manipulates Academia to Avoid Regulation. *The Intercept*. <https://theintercept.com/2019/12/20/mit-ethical-ai-artificial-intelligence/>
- O'Donnell, J., & Crownhart, C. (2025, May 20). We did the math on AI's energy footprint. Here's the story you haven't heard. *MIT Technology Review*. <https://www.technologyreview.com/2025/05/20/1116327/ai-energy-usage-climate-footprint-big-tech/>
- Pang, J., Ye, F., Wong, D.F., Yu, D., Shi, S., Tu, Z., & Wang, L. (2025). Salute the classic: Revisiting challenges of machine translation in the age of large language models. *Transactions of the Association for Computational Linguistics*, 13, 73–95. <https://aclanthology.org/2025.tacl-1.4/>
- Panwar, A. (2025, January 8). Generative AI and copyright issues globally: ANI Media v OpenAI. *Tech Policy Press*. <https://www.techpolicy.press/generative-ai-and-copyright-issues-globally-ani-media-v-openai/>
- Paulusse, M. (2026, February 24). Are AI-generated summaries suitable for studying and research? *Eindhoven University of Technology*, The Netherlands. <https://www.tue.nl/en/our-university/library/library-news/24-02-2026-are-ai-generated-summaries-suitable-for-studying-and-research>
- Pearl, J. (2018). *The book of why: The new science of cause and effect*. Basic Books.
- Peters, U., & Chin-Yee, B. (2025). Generalization bias in large language model summarization of scientific research. *Royal Society Open Science*, 12(4): 241776. <https://doi.org/10.1098/rsos.241776>
- Placani, A. (2024). Anthropomorphism in AI: Hype and fallacy. *AI Ethics*, 4, 691–698. <https://doi.org/10.1007/s43681-024-00419-4>
- Reuel, A., Ma, D. (2024). Fairness in reinforcement learning: A survey. In *Proceedings of the Seventh*

- AAAI/ACM Conference on AI, Ethics, and Society, AIES 2024, 1218–1230. <https://doi.org/10.5555/3716662.3716769>
- Russell, S.J., & Norvig, P. (2020). *Artificial intelligence: A modern approach*. Fourth edition. Pearson.
- Sadin, É. (2023). *Dissenting: Politics of our own* (Original: *Hacer disidencia: Una política de nosotros mismos*). Herder.
- Sadowski, J. (2025). *The mechanic and the luddite: A ruthless criticism of technology and capitalism*. University of California Press.
- Shah, C., & Bender, E. (2022). Situating search. In *Proceedings of the 2022 Conference on Human Information Interaction and Retrieval, CHIIR '22*, 221–232. <https://doi.org/10.1145/3498366.3505816>
- Shah, C., & Bender, E. (2024). Envisioning information access systems: What makes for good tools and a healthy web? *ACM Transactions on the Web*, 18(3), 1–24, 33, <https://doi.org/10.1145/3649468>
- Shojaee, P., Mirzadeh, I., Alizadeh, K., Horton, M., Bengio, S., & Farajtabar, M. (2025). The illusion of thinking: Understanding the strengths and limitations of reasoning models via the lens of problem complexity. *Appel Machine Learning Research*. <https://machinelearning.apple.com/research/illusion-of-thinking>
- Stadler, M., Bannert, M., & Sailer, M. (2024). Cognitive ease at a cost: LLMs reduce mental effort but compromise depth in student scientific inquiry. *Computers in Human Behavior*, 160, 108386, <https://doi.org/10.1016/j.chb.2024.108386>
- Tully, S.M., Longoni, C., & Appel, G. (2025). Lower artificial intelligence literacy predicts greater AI receptivity. *Journal of Marketing*, 89(5), 1–20. <https://doi.org/10.1177/00222429251314491>
- UNESCO (2022). Recommendation on the ethics of artificial intelligence. SHS/BIO/PI/2021/1, UNESCO. <https://unesdoc.unesco.org/ark:/48223/pf0000381137>
- Vassel, F.-M., Shieh, E., Sugimoto, C. R., & Monroe-White, T. (2024). The psychosocial impacts of generative AI harms. In *Proceedings of the AAAI Symposium Series*, 3(1), 440–447. <https://doi.org/10.1609/aaais.v3i1.31251>
- van Maanen, G. (2022). AI ethics, ethics washing, and the need to politicize data ethics. *Digital Society*, 1(2), 9. <https://doi.org/10.1007/s44206-022-00013-3>
- Verma, S., Boonsanong, V., Hoang, M., Hines, K., Dickerson, J., & Shah, C. (2024). Counterfactual explanations and algorithmic recourses for machine learning: A review. *ACM Computing Surveys*, 56(12), 312, 1–42. <https://doi.org/10.1145/367711>
- Walther, C.C. (2024, November 12). The hidden cost of AI energy consumption. *Knowledge at Wharton*. <https://knowledge.wharton.upenn.edu/article/the-hidden-cost-of-ai-energy-consumption/>
- Watermeyer, R., Phipps, L., Lanclos, D., & Knight, C. (2023). Generative AI and the automating of academia. *Postdigital Science and Education*, 6, 446–466. <https://doi.org/10.1007/s42438-023-00440-6>
- Watters, A. (2021). *Teaching machines: The history of personalized learning*. The MIT Press.
- Weidlich, J., Gašević, D., Drachsler, H., & Kirschner, P. (2025). ChatGPT in education: An effect in search of a cause. *Journal of Computer Assisted Learning*, 41(5), e70105. <https://doi.org/10.1111/jcal.70105>
- Weizenbaum, J. (1966). ELIZA—a computer program for the study of natural language communication between man and machine. *Communications of the ACM*, 9(1), 36–45. <https://web.stanford.edu/class/cs124/p36-weizenbaum.pdf>
- Weizenbaum, J. (1976). *Computer power and human reason: From judgment to calculation*. Penguin Books.
- Wieczorek, M., Hosseini, M., & Gordijn, B. (2025). Unpacking the ethics of using AI in primary and secondary education: a systematic literature review. *AI Ethics*, 5, 4693–4711. <https://doi.org/10.1007/s43681-025-00770-0>

- Wiener, N. (1960). Some moral and technical consequences of automation. *Science*, 131(3410), 1355–1358. <https://doi.org/10.1126/science.131.3410.1355>
- Wiggers, K. (2024, August 14). Study suggests that even the best AI models hallucinate a bunch. *TechCrunch*. <https://techcrunch.com/2024/08/14/study-suggests-that-even-the-best-ai-models-hallucinate-a-bunch/>
- Wilkins, J. (2025, October 19). The AI industry is traumatizing desperate contractors in the developing world for pennies. *Futurism*. <https://futurism.com/artificial-intelligence/ai-industry-traumatizing-contractors>
- Williamson, B., Molnar, A., Boninger, F. (2024, March 5). Time for a pause: Without effective public oversight, AI in schools will do more harm than good. *National Education Policy Center*. <https://nepc.colorado.edu/publication/ai>
- Winograd, T.A. & Flores, F. (1986). *Understanding computers and cognition: A new foundation for design*. Ablex Publishing Corporation.
- Winston, P. H. (1992). *Artificial intelligence*. Third edition. Addison-Wesley.
- Wong, M. (2025, March 7). Chatbots are academically dishonest. *The Atlantic*. <https://www.theatlantic.com/newsletters/archive/2025/03/chatbots-cheating-benchmark-tests-atlantic-intelligence/681954/>
- Wu, Z., Qiu, L., Ross, A., Akyürek, E., Chen, B., Wang, B., Kim, N., Andreas, J., Kim, Y. (2024). Reasoning or reciting? Exploring the capabilities and limitations of language models through counterfactual tasks. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, 1*, 1819–1862. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2024.naacl-long.102>
- Xiao, T., Nerini, F.F., Matthews, H.D., Tavoni, M. & You, P. (2025). Environmental impact and net-zero pathways for sustainable artificial intelligence servers in the USA. *Nature Sustainability*, 8, 1541–1553. <https://doi.org/10.1038/s41893-025-01681-y>
- Zewe, A. (2025, January 17). Explained: Generative AI's environmental impact. *MIT News*. <https://news.mit.edu/2025/explained-generative-ai-environmental-impact-0117>
- Zhao, C., Tan, Z., Ma, P., Li, D., Jiang, B., Wang, Y., Yang, Y., & Liu, H. (2025). Is chain-of-thought reasoning of LLMs a mirage? A data distribution lens. *Hugging Face*. <https://huggingface.co/papers/2508.01191>

## Acknowledgement

I am very grateful for the reviewers' and editor's comments on a previous version of this article.

## Author's Contributions (CRediT)

DM: Conceptualization, formal analysis, investigation, methodology, project administration, resources, validation, and writing, review, and editing. The author has read and agreed to the published version of the manuscript.

## Competing Interests

The author has no competing interests to declare.

## Acknowledgement of Use of Generative AI Tools

GenAI was not used in the development of this article.