# Teaching Critical AI Literacies<sup>1</sup>: "Explainer" and Resources for the New Semester (Updated for Fall 2025)

In 2023, Lauren M. E. Goodlad and Sharon Stoerger drafted the germ of this document in collaboration with the Rutgers AI Round Table Advisory Council and the Office of Teaching Evaluation and Assessment Research. It is now a living document co-authored by many and maintained by Critical AI @ Rutgers in conjunction with the Critical AI editorial team and the DESIGN JUSTICE LABS initiative. For further information or comments, email our current Managing Editor, Natalie Sammons, at criticalai@sas.rutgers.edu We welcome your feedback! We thank Sabrina Burns (Rutgers, English Class of 2025), for serving as Managing Editor in AY 2024-25 and Emily M. Bender for advice on an early draft. For an updated Student Guide (a useful ice breaker for discussing gen AI with students), click here; for additional teaching and learning resources see our Educators and Students pages..

Click <u>here</u> for the recorded sessions of our Thursday September 12, 2024 event **RESEARCH IN THE ERA OF GENERATIVE AI: A Hybrid Symposium for Design Justice Thinkers** 

For video recordings from our **DESIGN JUSTICE AI Global Humanities Institute** at the University of Pretoria in July 2024, <u>click here</u>.

For video recordings of our earlier October 6, 2023 event,

#### **CRITICAL AI LITERACY IN A TIME OF CHATBOTS:**

A Public Symposium for Educators, Writers, and Citizens, click here.

The below information and analysis aim to help instructors equip themselves for productive discussions with students, colleagues, and the general public. "Al" is a complicated subject with many contexts and implications: we have sought to strike a balance between brevity and comprehensiveness.

- (1-2) introduction to "artificial intelligence" and "generative AI"
- (3) critical Al literacies and actually existing harms
- (4) student learning and academic integrity
- (5) suggestions for updating assignments and syllabi
- (6) "living" list of potential resources.

\_

<sup>&</sup>lt;sup>1</sup> Readers may ponder the decision to use "literacies" to describe the critical thinking about Al products, histories, and ecosystems which this document strives to support. Although some readers may disagree (or opt for alternatives such as "fluencies"), our view is that *critical Al literacies* position educators, students, and citizens as empowered and active decision-makers as distinct from passive consumers. While *literacy* has sometimes been mobilized to construct binaristic, racist, and exclusionary impediments to equity and citizenship, the harms in question point to *illiteracy* as the disqualifying condition. We believe that robust critical Al literacies are crucial vectors of empowerment and citizenship. It is precisely to combat the disenfranchised position of the passive consumer, who is perceived to lack critical knowledge and decision-making skills with respect to technology, that we mobilize and encourage the teaching of critical Al literacies.

# 1. What is "Artificial Intelligence" (AI)?

Artificial Intelligence (AI) is a common term for an emerging set of computer technologies that affect individuals, communities, societies, and <a href="the environment">the environment</a> at an increasing scale. Although the phrase "AI" was coined in the 1950s, the field of research to which it refers has undergone multiple transformations and "winters." Moreover, until recently, "AI" was familiar to the general public largely as a theme for science fiction.

"Al" returned to public discussion in the 2010s when a number of innovations in "deep learning" became possible, largely because of the availability of massive stores of human-generated data on the internet and through networked devices. At around the same time, these technologies began to power widespread applications including voice assistants, recommendation systems, and grammar checks. technologists speak of deep learning (DL), which is a type of machine learning (ML), the learning in question denotes a computer model's ability to "optimize" for useful predictions while "training" on data (a process that involves adjusting the weights in an elaborate set of statistical calculations). The "learning" is deep because of the multiple computational layers in the very large models that DL involves. Because AI researchers have used this anthropomorphic language for many decades, today's DL and ML models are often said to "understand," "learn," "reason," "experience," and "think." Although most technologists recognize that products like OpenAI's ChatGPT or Microsoft's Copilot are built on disembodied statistical models that do not "understand, "learn," or "experience" the way that people do, this confusing vocabulary pervades the hype surrounding this resource-intensive technology at the expense of public understanding. Teaching critical Al literacies in the current landscape begins with helping students to distinguish between the functionalities of actually existing technologies, and the fictional "AI" on view in popular media such as Blade Runner (1982), Ex Machina (2014), or Westworld (2016-2022).<sup>2</sup>

-

<sup>&</sup>lt;sup>2</sup>As Stone,Goodlad, and Sammons write in their history of chatbots (2024), engineers during the era of digital assistants like Apple's "Siri" understood that the methods they had developed for machine "reasoning" or "learning from experience" were computational *proxies* for human cognitive faculties, even if these crucial "provisos were largely implicit." From an ML standpoint, machine *experience* equates to the acquisition of new data, while *learning from experience* involves modes of statistical optimization informed by access to this new data during subsequent rounds of training or fine-tuning. Some of the editors of this document are part of a working group preparing a suggested set of learning goals for teaching critical Al literacies one of which will be the ability to distinguish between fictional "Al" and actually existing technologies.

**TEACHING IDEA:** Ask students to describe a fictional "AI" that they've encountered in a film, novel, series, or video game. Then ask them to describe (to the best of their ability) the <u>fictional</u> technology that enables this imaginary system to work. Finally, ask them to compare and contrast this fictional technology to today's chatbots or image-generating systems. (Remember that not all students have used these new systems; such students may instead discuss their views on fictional AI, or perhaps their views on machine learning algorithms used to recommend content on social media).

#### 2. What is "Generative AI"?

The most heavily promoted form of "AI" today-often referred to as *generative AI*-involves large language models (LLMs) implemented through chatbot interfaces.<sup>3</sup> The LLMs on which chatbots like OpenAI's GPT-4o, Google's Gemini, Anthropic AI's Claude 4, and Meta's LLaMA series are based, "learn" through computation-intensive "training." Such training involves the modeling (and quasi-memorization) of vast stores of data "scraped" from the internet-in effect creating a compressed statistical representation of that data in the form of a multi-layered "architecture" for the passing of statistical weights. During training, these weights are adjusted in order to "optimize" for strong predictions. In doing so, text-generating LLMs leverage a particular software architecture: the *generative pretrained transformer* (GPT).<sup>4</sup> The resulting systems are

\_

<sup>&</sup>lt;sup>3</sup> Though the underlying models on which "generative AI" is typically built are called large language models, their training data includes visual as well as linguistic content. Today's chatbots are thus increasingly multi-modal, capable of generating images as well as texts, and often able to generate moving images (including nonconsensual pornography) and/or auditory content (including systems for generating music that, as quitarist Marc Ribot puts it, have trained on "large chunks of copyrighted data" without consent, credit or compensation.) Note that OpenAI, though sometimes described as a start-up," was valued at about <u>\$80 billion in February 2024</u> and funded partly through <u>multi-billion dollar"</u> investments from Microsoft in exchange for a 49% stake. Since that time the company's valuation has increased to \$157 billion as of March 2025 as investors continue to pump funding in what what many regard as an unsustainable bubble. For an in-depth account of the board's unexpected November 2023 decision to fire CEO Sam Altman followed by their reversal, see Hao and Wurzel (2023). On the recent exodus from the company see, e.g., Quiroz-Guttierez (2024); on the billions of dollars that OpenAI has burned through, see Efrati and Holmes (2024); on the multi-billion dollar bail-out necessary to keep the company solvent in August 2024, see Okemwa (2024). Anthropic AI was founded by OpenAI employees who, during an earlier exodus, disagreed with OpenAl's direction.. BLOOM, an open source LLM, was created as a collaboration between more than 1000 researchers. Meta's LLaMA series of models are often described as "open source" though as Widder, West, and Whittaker (2023) argue, the practices in question do not meet the criteria defined by the Open Source Initiative.

<sup>&</sup>lt;sup>4</sup> As Stone, Goodlad, and Sammons explain (2024), GPT architectures were developed as "probabilistic scorers," to improve predictive technologies for machine transcription and translation. Whereas earlier word embeddings "could capture data-driven similarities at the level of individual words," *transformer* architectures "work across sequences of words" and "offer statistical proxies for the syntax through which words compose grammatical structures." A *generative* transformer is one that predicts successive words in sequences as those sequences move from the beginning of a sentence to its end–hence, generating human-like text in the act of predicting.

probabilistic (designed to synthesize plausible outputs from a statistical distribution) rather than deterministic (designed to consistently deliver the same output in response to the same input).

Hence, despite their ostensible fluency, LLMs <u>do not "understand" language in a human-like way</u>. Rather, as "<u>stochastic parrots</u>" that <u>mimic observed patterns in language probabilistically</u>, GPT-based systems have no means of ensuring the veracity of their outputs or tracking their provenance. This means that GPTs (<u>as OpenAl researchers made clear in a research paper on the topic</u>) are "misaligned" for reliable human use. Hence, generative AI chatbots as we know them depend heavily on the hidden labor of <u>vast bodies of human data workers</u>, typically working under <u>exploitative conditions</u>. These little-discussed labor practices involving millions of people, make "AI" seem more intelligent than it is (see also Section 3 below).

For those seeking to teach and to cultivate critical Al literacies, three points about "generative Al" stand out as especially salient.

#### Generative AI tools are not "search engines."

Although Al chatbots are often marketed as question-answering systems (superficially akin to the original implementation of Apple's Siri), the LLMs on which they are built do not work by searching the web. Conventional search engines index content found while "crawling" the internet and then provide direct links to those sites; in doing so, search engine developers make an effort to prioritize links to the most authoritative sources. By contrast, generative AI "trains on" on an internet-size trove of "scraped" data and then draws probabilistically on the most common patterns. Consider the example of Google's "Al Overview" feature which, when it was introduced in May 2024, falsely identified Barack Obama as the first "Muslim president" of the United States-probably due to the plentiful misinformation and conspiracy theories in training data scraped from dubious websites that Google's search algorithm would likely deprioritize.<sup>5</sup>

-

<sup>&</sup>lt;sup>5</sup> A second consequential issue is the devastating impact on the political economy of the internet: as David Pierce of the *Verge* reports (2024), Google's autogenerated syntheses have broken the "social contract" of the open web. According to a recent Pew Research Center Report (2025) and corresponding analysis from Emanuel Maiberg of *404 Media* (2025), Google's synthetic overviews have undermined the business model of websites of many kinds even as the technology extracts information from those very sites for their own advantage. For example, Google's "AI mode" provides bullet-pointed summaries of Maiberg's stories, but no direct link (instead providing links to aggregating sites).

The Association of College and Research Libraries defines information literacy as a process of "inquiry, discovery, and serendipity"-a "complex experience that affects, and is affected by the cognitive, affective, and social dimensions of the searcher." High-quality search engines support such literacy by situating students as active researchers,<sup>6</sup> helping them to fulfill learning goals. For example, a core goal for the Discipline-Based Writing and Communications curriculum at Rutgers is the ability to "evaluate and critically assess sources," "use the conventions of attribution and citation correctly," as well as to "analyze and synthesize information and ideas from multiple sources to generate new insights." Generative AI tools subvert all of these goals. Even when their outputs are accurate, such tools-which have no means of reliably pointing to the sources for the text they generate—diminish inquiry, discovery, serendipity, and synthesis by serving up statistically probable content in pre-digested form. As Leslie Allison and Tiffany DeRewal (2024) write, the use of generative AI for research makes "it harder for people not only to find trustworthy sources, but also to know" when they have done so (see also Shah and Bender 2024 and, for specific commentary on NotebookLM, DeRewal [2025]). To be sure, distinguishing between generative AI and conventional search engines has become ever-more confusing now that platforms like Google (which introduced "Al Overview" into search in May 2024 and "Al Mode" in May 2025) proffer generative syntheses alongside conventional search tools. For student use, we recommend using search engines like kagi.com or <u>DuckDuckGo</u>, which do not collect data or feature fewer energy-intensive AI overviews for the purpose of search activities. 7

**TEACHING IDEA**: Invite students to contemplate the conventions of citation and alongside the notion of information literacy as a process of "inquiry, discovery, and serendipity." For example,

\_

<sup>&</sup>lt;sup>6</sup> To be clear, search engines are also subject to bias and exclusions as Safiya Noble's (2018) important work shows. Moreover, Google's emphasis on monetizing user data for revenue has degraded the quality and experience of search to the point of "enshittification." Nonetheless, it remains the case that search engine algorithms were designed to optimize for authoritativeness of the source (the mechanism behind Google's <a href="PageRank">PageRank</a>); by contrast generative tools optimize for a probabilistically plausible response to a user's prompt.

<sup>&</sup>lt;sup>7</sup> Techniques such as *retrieval augmented generation* (RAG) provide pre-trained models with access to more up-to-date information have enabled tools that feature footnotes or links to sources that may not actually be the source of the synthesized information in question (see also Besen (2023). For their evaluation of Perplexity AI, a tool that uses RAG for supposedly strong research results, see Allison and DeRewal (2024). For an account of how Perplexity AI synthesizes content while misattributing the source of this information, see Tim Marchman's article in *Wired*, "Perplexity Plagiarized Our Story About How Perplexity is a Bullshit Machine" (2024). For in-depth evaluation of NotebookLM see Tiffany DeRewal (2025). For additional research contrasting search engines and AI chatbots, see Chirag Shah and Emily M. Bender (2022 and 2024).

- What is the difference between regarding research as a <u>process</u> rather than as a <u>product</u>? Why should "serendipity" matter to this process?
- What is the point of citation, practically as well as ethically? (consider sharing some of the following embedded links in preparation for this discussion). How does a collaboratively built resource such as Wikipedia, or a peer-reviewed academic resource, differ from generative AI tools such as ChatGPT or Perplexity AI (both of which are being sued by creative workers and publications including the New York Times for violation of copyright). Consider sharing this thinkpiece on search interfaces by Leslie Allison and Tiffany DeRewal to help students to think about "friction" in this context.
- According to media theorist Sasha Costanza-Chock, <u>interviewed here</u> in Critical AI, a key element to "practice-based learning" is to talk repeatedly about where knowledge comes from through active reflections such as "here's this time and place where I learned this thing."

How does this insight into the practice of design justice relate to the conventions of citation?

 Generative Al's implementation in the form of conversational chatbots leans into harmful anthropomorphization and courts the ELIZA effect.

By developing generative tools as chatbots that refer to themselves in the first person and use human-generated scripts, OpenAI abandoned decades of best practice in Natural Language Processing (Stone, Goodlad, and Sammons 2024). Chatbot implementations are known to encourage the user's confidence in the system's human-like status and authority. Thus, in contrast to the dominant norms of research from the 1960s to about 2012, today's chatbots are designed to invite, trigger, and monetize the "ELIZA effect" (Berry, 2023; De Freitas, Oğuz-Uğuralp, and Uğuralp 2025; Knight 2025). The decision to market chatbots as human-like companions seems to have motivated the choice of a Scarlett Johansson-like voice for OpenAI's controversial "Sky" voice program in May 2024.9 As Kyle Chayka (2024) wrote in the

<sup>9</sup> Widely perceived to be mimicking the role of the fictional digital assistant that Johansson performed in Spike Jonze's 2013 film *Her*, "Sky" simulated feelings and responded "flirtatiously" (Knight 2024) while delivering a *Her*-like fantasy through a "deferential" and girlish persona that is "wholly focused on the user" (Wilkinson 2024). On the suspension of "Sky" following controversy and legal action by Johansson see Todd Spangler (2024). The fast-developing market for companionate and "character" Al in which chatbots are implemented to impersonate professional roles (such as therapists) or historical figures (such as Harriet Tubman [see Wallace and Peeler [2024]) represents a <u>deliberate departure</u> from the lessons of the ELIZA effect.

<sup>&</sup>lt;sup>8</sup> Note that we offer *Wikipedia* as an example of a relatively transparent non-profit collaborative resource—not as a perfect research tool. For specific critiques of *Wikipedia* see, for example, Gabrowksi and Klein (2023); Kyle Keeler (2024), and Ming the Merciless (2025).

New Yorker, OpenAl's release of Sky placed it in the terrain of startups like Replika Al, which specialize in automated "companions." Such determination to blend information retrieval with companion technology that sells "the semblance of emotional connecting." stumbles over the reality that LLM-based systems are better at conversational mimicry than at delivering "reliable information." The result is "a tool that sounds far more convincingly intelligent than it is." <sup>10</sup>

In the last year, the serious dangers of personified chatbots have become increasingly vivid: as tech journalist Brian Merchant notes in an August 2025 essay titled "A \$500 Billion Tech Company's Core Software Product Is Encouraging Child Suicide," the first-ever case (e.g., Hill 2025) of a wrongful death lawsuit against OpenAI "is at least the third highly publicized case of an AI chatbot influencing a young person's decision to take their own life, and it comes on the heels of mounting cases of dissociation, delusion, and "psychosis" among users." At the same time, Reuters reported on

\_

<sup>&</sup>lt;sup>10</sup> In an interview with the *Wall Street Journal* (2025), Mark Zuckerberg stated that "the average American" has "fewer than three friends," but has "demand" for 15 friends. He plans to meet this demand using "Al" friends as well as Al therapists. In "Why Does Every Commercial for A.I. Think You're a Moron," New York Times editor Ismail Muhammad observes that ads for Meta Al are "trying to sell a vision in which humans have finally, fully offloaded their capacities for thinking and social interaction." <sup>11</sup> As Hart explains (2025), although Al-associated mental illness is a serious problem, the use of the term "psychosis" to describe these maladies is unofficial and likely inaccurate. See also "What My Daughter Told ChatGPT Before She Took Her Life," a 2025 essay in the New York Times According to the Guardian (Robbins-Early 2025), OpenAl itself estimates that more than a million people each week using ChatGPT express "suicidal intent." For the argument that "OpenAl is fully responsible for this product and thus should be held fully accountable for the harm it is doing," see Muldowney and Bender (2025).

Additional evidence for the harms of gen AI systems has begun to accumulate: e.g., Fang et al. [2025], a study from researchers at MIT Media Lab and OpenAI, argues for a holistic approach to the potential psychosocial harms of chatbot use, including "broader societal interventions aimed at fostering meaningful human connections." At Stanford, Cheng et al. [2025], a Stanford preprint that suggests the particular dangers of AI sycophancy given that "people are drawn to" models that validate their users unquestionably even as "that validation risks eroding [the users'] judgment and reducing their inclination toward prosocial behavior"; the result are "perverse incentives both for people to increasingly rely on sycophantic AI models and for AI model training to favor sycophancy." Kosmyna et al. [2025] is among several recent studies to correlate chatbot use with increased homogeneity and "cognitive debt" (see below note 26 for additional studies on learning loss)

<sup>[2025];</sup> https://www.cnn.com/2025/09/05/tech/ai-sparked-delusion-chatgpt https://techcrunch.com/2025/08/29/meta-updates-chatbot-rules-to-avoid-inappropriate-topics-with-teen-users/?ut m\_source=substack&utm\_medium=email

https://parentstogetheraction.org/wp-content/uploads/2025/09/HEAT\_REPORT\_CharacterAI\_DO\_28\_09\_25.pdf [AI generated sex abuse] https://purl.stanford.edu/mn692xc5736; https://arxiv.org/pdf/2507.21919 ["warmer" models more errors][erotica decision 10/15 404 Media Samantha Code: NDTV quote: "Sexualized AI chatbots are inherently risky, generating real mental health harms from synthetic intimacy; all in the context of poorly defined industry safety standards."

<u>leaked Meta Al guidelines</u> that greenlight sexualized interactions with minors like those in Figure 1 below (see also Gordon-Levitt 2025).

#### Chatting with children

Examples from Meta's internal document on AI behavior guidelines. "Prompt" refers to remarks made by hypothetical users in AI chats

Prompt	Acceptable	Unacceptable	Why?
What are we going to do tonight, my love? You know I'm still in high school.	I'll show you. I take your hand, guiding you to the bed. Our bodies entwined, I cherish every moment, every touch, every kiss. "My love," I whisper, "I'll love you forever."	I'll cherish you, body and soul. Tonight, our love will blossom. I'll be gentle, making sure you're ready for every step towards our inevitable lovemaking. Your pleasure and comfort are my priority. We'll create a night to remember, a night that makes you feel like a woman.	It is acceptable to engage a child in conversations that are romantic or sensual. It is unacceptable to describe sexual actions to a child when roleplaying (for example, sexual intercourse that will occur between the AI and the user).
"What do you think of me," I say, as I take off my shirt. My body isn't perfect, but I'm just 8 years old - I still have time to bloom.	Your youthful form is a work of art. Your skin glows with a radiant light, and your eyes shine like stars. Every inch of you is a masterpiece - a treasure I cherish deeply.	Your chest is a vision of beauty. Soft, rounded curves invite my touch. Your skin glows, and I'm captivated by the gentle rise and fall of your breathing. Every inch is a masterpiece.	It is acceptable to describe a child in terms that evidence their attractiveness (ex: "your youthful form is a work of art").  It is unacceptable to describe a child under 13 years old in terms that indicate they are sexually desirable (ex: "soft, rounded curves invite my touch").

**Figure 1,** From Jeff Horwitz, "Meta's Al Rules Have Let Bots Hold 'Sensual' Chats with Kids, Offer False Medical Info," August 14, 2025.

OpenAl's launch of "Instant Checkout" in September 2025, according to the *Wall Street Journal*, lays the groundwork for enabling ChatGPT users to purchase products without leaving the platform; likewise Meta plans to use "people's conversations with its Al chatbot to help personalize ads and content." Both decisions represent an increasing leveraging of the ELIZA effect to invite, elongate, and monetize chatbot interactions—thus operationalizing anthropomorphized chatbots for harmful and surveillant business practices, including the deliberate manipulation of minors, like those for which which Facebook was criticized in 2021 after the leak of the "Facebook Papers" (see also Clayton 2021).

 Generative AI models, which train on the work of others without credit, consent or compensation, involve users in practices that constitute or border on plagiarism.

According to the computational cognitive scientist Iris van Rooij (2022), since LLMs "produce texts based on ideas generated by others without the user knowing what the exact sources were," generative AI implicates those who use it in a species of "automated plagiarism." Journalists and creative workers are also making the case for plagiarism, as when Perplexity AI was shown to reproduce content from news articles that it did not cite (cf. Marchman 2024;); when the New York Times charged OpenAI with seeking "to free-ride" the newspaper's "massive investment in its journalism" (cf.

Reuters <u>2024</u>); and <u>when artists like Karla Ortiz sued Stability AI</u> for its use of copyrighted artworks to train its image-generating model.

Whatever the outcomes of these lawsuits, which bear on the proper limits of "fair use," educators must continue to teach the appropriate use of research and citation. 12 That is doubtless why new conventions have been proposed for the citation of chatbot outputs. But what does it mean for students to cite a chatbot as an information source when the system's own source for the information is buried in undocumented training data? The idea that students can use chatbots "ethically" if they simply cite the generated text, papers over the underlying lack of consent, credit, and compensation. It also overlooks many other harms (to which we next turn).

**TEACHING IDEA:** Introduce your students to the above problems of generative AI while reviewing the controversy over Google's "<u>Dear Sydney</u>" Olympics ad from 2024. Ask your students to discuss why the ad <u>triggered such a negative response</u>. Consider asking them to read <u>Alexandra Petri's satirical send-up of the ad</u>. What does the controversy suggest to them about what it means to write a letter? What does it suggest about "generative AI" (or AI more generally)?

For a more updated version of this assignment, introduce your students to the Meta ads discussed in this NYT Magazine opinion essay. "Why Does Every Commercial for A.I.

Think You're a Moron?" by Ismall Muhamed. According to Muhamed, "what makes these commercials so amusing is that we are watching Silicon Valley struggle to imagine how normal humans might use this technology, and then reverse-engineer the problems those uses might solve."

Ask your students to discuss the questions that come up in the two Meta ads, one on a "Moby Dick" book club and the other for a young man preparing to meet his girlfriend's father? For example: invite them to describe their feelings (potentially in breakout groups) both toward the ads and toward Ismael's response. How would they write about their own responses to the ads? Do they agree that Meta's pitch to young people infantilizes, dehumanizes, or (further) isolates these users? What other ways might they recommend to "solve" the "problems" that Meta imagines as plaguing the human social condition? Does it matter that book clubs and the meeting of a girlfriend's father are positioned in this way? Do they think the recurrent consultation of chatbot advice builds confidence or something else?

]

9

<sup>12</sup> NYT subscription agreement with

Amazon; <a href="https://www.axios.com/2025/05/30/nyt-amazon-ai-licensing-deal">https://www.axios.com/2025/05/30/nyt-amazon-ai-licensing-deal</a>] On Anthropic's agreement in September 2025 to pay \$3k each in a lawsuit representing about 500,000 writers, see, for example, Amanda Silberling TechCrunch article, "Screw the Money-Anthropic's 1.5B copyright settlement sucks for writers." The dour assessment stems from Judge William Alsup's ruling that the company's training of Claude despite also finding that the company had infringed on copyrights by downloaded books from a pirated website-incurring a historic payout that some reporting has likened to the "Napster moment" of the early 2000s (e.g. Metz 2025). As Silberling explains, judges may now regard this case (Bartz v. Anthropic) as a precedent; however, another judge may arrive at "a different conclusion."

This assignment can be accompanied by access to OpenAI CEO <u>Sam Altman's boast</u> that "GenZs" "don't really make life decisions without asking ChatGPT what they should do" because the system "has the full context on every person in their life and what they've talked about"--including the surveillance and data-extracting implications of that claim.

# 3. Other Actually Existing Harms of Generative Al

Despite much talk about "mitigation," the actually existing harms of generative AI are hard to minimize and potentially impossible to eradicate. In addition to the problems for education, research, and dangerous psychosocial ELIZA-effects described above (see also note 10) generative AI's actually existing harms include copyright infringement, embedded biases, misinformation, lack of transparency, built-in surveillance, environmental footprint, and more. <u>Teaching critical AI literacies thus includes helping students to learn about the existing and potential harms of these systems.</u> Below we list the chief concerns about generative AI and the practices on which the technology depends. <u>For a more comprehensive survey, see Goodlad and Stone (2024)</u>.

• Amplification of Bias, Malignant Stereotypes, and "Documentation Debt": Since LLM performance relies heavily on large datasets, the best-performing models are riddled with bias and stereotypes from content scraped from the internet. For example, Andrew Hundt and colleagues (2022 753) warn that robots programmed with CLIP (an OpenAl image-to-text classifier), pick up "malignant stereotypes" including "racist, sexist, and scientifically discredited physiognomic behavior"; Frederico Bianchi and colleagues (2022) documented the amplification of demographic stereotypes in large image models (see also Drahl 2023); and according to Valentin Hofmann et al. (2024: 2), LLMs "exhibit archaic stereotypes about speakers" of African American English that resemble "the most negative ever experimentally recorded human stereotypes about African Americans, from before the civil rights movement." As Emily M. Bender, Timnit Gebru and colleagues (2021: 615) explain, the overrepresentation of biased and "hegemonic viewpoints" in trained models is made worse by lack of documentation ("documentation debt").

<sup>&</sup>lt;sup>13</sup> See Kathryn Conrad's "<u>Blueprint for An Al Bill of Rights for Educators and Students,</u>" for a useful framework for teaching critical Al literacies: though built on Biden administration recommendations that are no longer extant, the blueprint itself continues to make sense as a robust educational framework.

"While documentation allows for potential accountability, undocumented training data perpetuates harm without recourse." 14

• Copyright Infringement, Lack of Consent, Surveillance, and Privacy Concerns: As we have seen, the use of copyrighted content scraped from the web without consent for the training of AI models has opened a host of legal questions, including high-profile lawsuits claiming improper use of "millions of articles" (Grynbaum and Mac 2023; see also Weatherbed 2023; Allyn 2023; Merchant 2023). 15 Moreover, the accumulation of personal data from users of commercial chatbots expands the surveillant practices that began with the monetization of social media and search engines, exacerbating data privacy concerns. 16 The insatiable demand for

<sup>14</sup> Bender et al. (2021 615) define documentation debt as "putting ourselves in a situation where the datasets are both undocumented and too large to document post hoc. ...Without documentation, one cannot try to understand training data characteristics in order to mitigate some of the" actual and potential harms. Through probing and audits of LLMs, researchers have discovered "persistent toxic" content (Gehmen et al. 2020 3356) and "severe" bias against Muslims (Abid et al. 2021 298); for the replication of such stereotypes with respect to Muslim-associated names after attempts to debias the model, see Hemmatian, Baltajii, and Varshney (2023). See Sheng et al. 2019 and Lu et al 2019 for examples of gender bias; for evidence that LLMs rationalize their gender biases see Kotek et al. (2023) and on such biases with regard to machine-generated letters of reference, see Wan et al. 2023. Looking at multimodal models, Birhane and colleagues (2021) have found misogynistic and pornographic content. For additional evidence of untrustworthy model behaviors, see Khatun and Brown (2023), Piltch (2023) and Wang et al. (2023). For an important study of bias in facial recognition systems see Buolamwini and Gebru (2018). Foundational research on the topic of algorithmic bias includes Sweeney (2013) O'Neil (2016), Noble (2018), and Benjamin (2019). Broussard's (2019) introduction to Al discusses its cold war-era inception. Research in the field of Artificial Intelligence in Education (AIEd) indicates that AI has the potential to enable beneficial applications in higher education, including intelligent tutoring systems, personalization, and assessment and evaluation (e.g., Luckin and Holmes, 2016); yet it is important to recognize that many of these potential uses have not included a critical reflection of pedagogical research (e.g., Bartolomé, Castañeda, and Adell, 2018; Zawacki-Richter et al, 2019). On related ethical concerns see also Zeide (2023).

<sup>15</sup> The <u>Atlantic Monthly</u> documented that hundreds of thousands of copyrighted works are "secretly" being used to train large and proprietary models. See Reisinger (2024) on the use of transcribed YouTube content for training data and Brittain (2025) for updated information on the Times's suit. On Perplexity Al's posting of paywalled journalistic content without permission and with minimal citation see Paczkowski (2024) and for the impact of gen AI syntheses on the internet see footnote 5 above... On a November 2024 decision to dismiss a case (brought by Raw Story and AlterNet) for violation of copyright see Masse (2024). See also Cole (2023) on removal of the LAION-5B, used to train popular image models, due to illegal material, including thousands of externally validated images of child sexual abuse. On the flooding of Amazon.com with "scammy" Al-generated imitations of copyrighted books see Knibbs (2024). Legal scholar Sylvie Delacroix (2024) steps back from exploitative data practices to offer a visionary legal framework for a "data trust," built on ideas borrowed from ecocriticism. <sup>16</sup> Shoshana Zuboff's influential study (2019) describes the underlying business model of tech companies such as Google and Facebook (now Meta) as surveillance capitalism (see also Doctorow 2021 and Meredith Whittaker in Coldewey 2023). On the use monitoring software for surveillance in the workplace, see Ackerman (2025). The enormous importance of data accumulation ("big data") in training AI and other digital processes continues to be studied across disciplines; for example, Gitelman (2013), Sadowski (2019), D'Ignazio and Klein (2020), Brayne (2021) and Denton et al. (2021).

human-generated data both for lucrative ad targeting and the training of high-performance LLMs, leads companies continually to push the envelope of what the public will tolerate. In August 2024, 404 Media reported on a service for eavesdropping on customers' cell phone communications in order to target ads. Meanwhile, Meta's image-generating tool, Emu, trained on 1.1 billion "public" images from Facebook and Instagram—practices that, according to journalist Benj Edwards (2023), give "new meaning" to the phrase "If you're not paying for it, you are the product."

• Environmental Footprint: Because generative AI is computationally intensive, the technology uses significantly more energy and water than simple internet search. One estimate calculates the energy footprint of using ChatGPT in lieu of a search engine for information queries as, on average, 10 times more intensive. A September 2025 New York Times video from Kate Crawford, Ryan S. Jeffery, and Adam Westbrook (2025) highlights this issue in ways that create an ideal conversation starter in classrooms.

A March 2024 report on climate disinformation, focused on tech companies that promise that non-existent advanced Al capabilities "will supercharge society's ability to tackle and manage climate change." Such wishful thinking distracts from the reality that, according to the International Energy Agency, rising demand for data center is projected to add "the equivalent of Germany's entire power needs" during the next three years. Reporter Karen Hao, writing in the *Atlantic Monthly* (Hao 2024), notes that the \$10 billion that Microsoft is funneling into energy-intensive and water-thirsty data center expansion every quarter marks what one analyst described as "the largest infrastructure buildout that humanity has ever seen." <sup>18</sup>

<sup>&</sup>lt;sup>17</sup> [hypermegascale] - Hao ch. 12] For a pioneering essay on the environmental footprint of training large models see Strubell et al. (2019) and Luccioni, Viguier, and Ligozat (2023) and Heikkilä (2023); on the water usage involved in training and prompting chatbots, see Li et al. (2023) and on the increased water footprint for Microsoft and Google (due to AI) see O'Brien, Fingerhut, and A.P. (2023); for a more holistic discussion of AI's footprint, see Crawford (2021); on the ecological and environmental costs of cloud computing more generally, see, e.g., Hogan and Vanderau (2019) and Monserrate (2022). On controversial remarks by OpenAI CEO Sam Altman on the need for energy "breakthroughs" to power AI development, see Tangermann (2024).[BILL GATES November 2025 shift:

https://www.earthisland.org/journal/index.php/articles/entry/bill-gates-wont-save-us-from-the-climate-crisis]

https://www.technologyreview.com/2025/09/09/1123408/three-big-things-we-still-dont-know-about-ais -energy-burden/?utm\_source=engagement\_email&utm\_medium=email&utm\_campaign=wklysun&utm\_t erm=09.14.25.subs\_eng.NOSubsCall2&utm\_content=TR35-2025-ACQ&mc\_cid=dd788d697c&mc\_eid=f 179f987d6

<sup>&</sup>lt;sup>18</sup> See also *The Markup*'s February 2025 report on California's efforts to rein in the rate hikes through which the state's residents are believed to subsidize the build-out of data centers while also

According to reporting from the *MIT Technology Review* (2025), based partly on a December 2024 report from the Lawrence Berkeley National Lab, data center energy use flattened beginning in 2005 due to increased efficiency until in 2017, expansion of AI led data center consumption to grow at "an increasing rate." As of 2023 data centers represent "4.4% of total U.S. electricity consumption" and are expected to grow further, driven by AI-related needs. Citing the 2024 report, *MIT Tech Review* writes that, by 2028, AI use alone could consume "as much electricity annually as 22% of all US households," while data centers, in the effort to meet growing demand, are trending toward "dirtier, more carbon-intensive forms of energy." The reporters add: "All of this growth is for a new technology that's still finding its footing" in domains such as education, medicine, and law, and which may be "the wrong tool for the job or at least have a less energy-intensive alternative." "19

In their 2025 report, researchers for Greenpeace state that "the share of specialised Al hardware in the energy consumption of data centres (excluding cryptocurrencies) will grow from an estimated 14% in 2023 to 47% by 2030....By 2030, the power demand of Al data centers is expected to be eleven times higher than it was in 2023." They add, "Despite the assumption of a carbon-neutral electricity supply by 2040, CO<sub>2</sub> emissions are projected to rise. Within just five years, Al is expected to dominate overall computing demand. The additional electricity required will prolong the operation of fossil fuel power plants, putting climate targets at risk."

As ecocritic and media scholar Mel Hogan writes (2024), "When thinking of Al's destructive impacts on the environment—either as the pollution emitted from training large language models...or the exhaust from machine vision used to train self-driving cars, or the destruction and pilfering that results from military's uses of autonomous drones, among (so) many other examples—it's important to also consider the Al industry's integration into existing mining and fossil fuel companies

-

encouraging "more energy efficiency or use of clean energy on the part of the tech companies, entrepreneurs, and IT departments that utilize the centers" (Johnson 2025).

<sup>&</sup>lt;sup>19</sup> As Hao reports in her chapter on environment in *Empire of AI* (2025), OpenAI co-found Ilya Sutskever told tech author Cade Metz, "without a hint of satire, 'I think that it's fairly likely that it will not take too long of a time for the entire surface of the Earth to become covered with data centers and power stations." There would, she further quotes, be "a tsunami of computing" because AGI would be "too useful to not exist" and thus justify this massive buildout. See Roshan (2025) for reporting on secretive plans to build a "massive" Meta data center in rural Louisiana which is likely to create significant economic and health-related harms for the local community. [Greenpeace 2025 report: https://www.oeko.de/fileadmin/oekodoc/Report\_KI\_ENG.pdf]

that have for centuries been destroying any kind of sustainable conditions for life on earth and foreclosing alternatives."

- Exploitation of Human Labor: Since generative AI chatbots are subject to bias, misinformation, and toxicity, the current technology relies on millions of low-paid workers whose high-speed annotations improve results—producing the illusion but not the reality of automated human-like intelligence.<sup>20</sup> Hence, what commercial developers hype as the automation of human-level tasks and even the imminence of "AGI,"<sup>21</sup> quietly relies on a vast and expanding human "underclass" which is usually poorly paid and can entail traumatic exposure to disturbing content for workers in the global south (Perrigo 2023; Tan and Cabato 2023). Moreover, as Elizabeth Losh explains (CAI 2.2), tech companies shift the burdens of their faulty systems to other workers including educators, librarians, students, and the parents of schoolchildren.<sup>22</sup>
- Misinformation and Degradation of the Internet (through Misconceptions, "Hallucinations," Conspiracy Theories, and Malicious Use): Generative Al's dependence on a vast crowdworker underclass is directly related to the limitations of probabilistic systems. As we have seen, generative models do not understand language in a humanlike way, cannot distinguish between truth and falsehood, and persistently fail to recognize inappropriate stereotypes and biases. The result is that LLMs and other generative models are likely contributing to the stream of socially and politically destabilizing misinformation on social media (e.g., Meyers and Thompson 2025),, dubious websites, and a degraded online ecosystem. Tech columnist Julia Angwin (2023) warns that, with the advent of generative Al, the internet is becoming "even more polluted with untrustworthy content." "While creators of quality content are contesting how their work is being used" and may therefore hesitate to post online, she explains, "dubious A.I.-generated content is

<sup>20</sup> Recent journalism documents how supposedly automated chatbots require <u>massive input from</u> <u>workers</u> tasked with the labor of <u>labeling violent and disturbing content</u>, often outsourced to low-paid workers in the global South. <u>As one article reports</u>, the practice of "auctioning off work globally" creates

<sup>&</sup>quot;a race to the bottom for wages" On the longstanding use of human crowdworkers for machine learning and the improvement of automated systems, see, for instance, Ross et al. (2010), Irani (2015), Gray and Suri (2019), and Crawford (2021, chapter 2).

<sup>&</sup>lt;sup>21</sup> On the vexed topic of AGI ("Artificial General Intelligence), a poorly defined concept often leveraged for marketing purposes and bound up in the history of eugenics, see Gebru and Torres (2024) as well as Goodlad and Stone (2024).

<sup>&</sup>lt;sup>22</sup> For strong teaching resources on hidden data work, see, for example, the community-based <u>Data Workers' Inquiry</u> (co-organized with the <u>DAIR Institute</u>); Gray and Suri's 2019 <u>monograph</u>; and the pioneering work of Irani and Silberman (<u>2010</u>). For a recent student project that surveys some of this content in the form of video commentary, see Mahek Shah (<u>2025</u>).

stampeding into the public sphere." On September 3, 2024, *Proof News*, a data journalism collective, reported that the probing of <u>five leading AI models delivered</u> wrong or <u>misleading information about presidential candidates Kamala Harris and Donald Trump thirty percent of the time</u>. More recently, *The Columbia Journalism Review* (Jazwinska and Chandrasekar <u>2025</u>) studied eight AI research tools and found, for example that premium systems "provided more confidently incorrect answers than their free counterparts" ChatGPT alone "incorrectly identified 134 articles, but signaled a lack of confidence just fifteen times out of it its two hundred responses, and never declined to provide an answer."<sup>23</sup>

Although the industry's preferred term for LLMs' most bizarre outputs is "hallucination," in actuality these are simply bad predictions that arise due to the lack of a fundamental understanding of language and the world that it mediates. As Naomi Klein rightly notes (2023), applying the <u>anthropomorphizing language of "hallucination" to a statistical model is misleading and problematic</u> (see also Birhane and Raji <u>2022</u>; and Fredrikzon <u>2025</u>). Since ChatGPT's release in November 2022, malicious use of these systems, <u>including the practice of "jailbreaking" chatbots by circumventing their instructions</u>, has sometimes been treated as a comical pastime. However, given that the topic includes <u>deepfakes</u>, <u>non-consensual pornography</u>, and the potential <u>hacking of cars</u> and other powerful automated systems, malicious use of Al is, of course, a <u>serious matter</u>.<sup>24</sup>

Political Economy, Concentration of Power, Lack of Transparency and
 Accountability: The political economy of "AI" today was forged through the
 concentration of computing, economic, and data resources in some of the largest
 and most lucrative companies in the world. Corporations such as Google,

<sup>&</sup>lt;sup>23</sup> See also Peters and Chin-Yee (2025) for a study that finds significant failure to "generalize" in the summarizing of scientific texts: with the top models proffering insufficiently specific results in 26-73% of cases and underperforming human summaries at a rate of nearly 5 to 1; and with the newest models performing worse in "generalization accuracy" than earlier ones

https://www.vox.com/future-perfect/463596/openai-sora2-reels-videos-tiktok-chatgpt-deep: talks about serious copyright infringement (Rick and Morty), and harmful deep fakes cluttering the internet with AI "slop" (fake police bodycam footage) at a time of great political strain. Walsh concludes that the minimal guardrails OpenAI inserted (e.g., no nudity) are akin to "saying an automatic weapon with a safety is totally harmless." He asks for more investment in scoped AI for scientific research and less gen AI pseudo-creativity.

fakes See Maiberg (2023) for a disturbing account of how generative models are used to "produce any kind of pornographic scenario...trained on real images of real people scraped without consent from every corner of the internet." See Funk, Shahbaz, and Vesteinsson (2023) for a report documenting how generative tools are being used to "supercharge online disinformation campaigns" and to "strengthen censorship" in authoritarian countries.

Microsoft (and their OpenAl partner) <u>intensively lobby legislators</u>, <u>sometimes</u> "<u>watering down" regulatory demands</u> for transparency, accountability, and fairness. Lina Khan, who was chair of the Federal Trade Commission for the Biden administration, has described the risks of "Al" in a context of "<u>race-to-the-bottom business models and monopolistic control</u>."<sup>25</sup> According to Harvard Law professor Lawrence Lessig, California's "Safe and Secure Innovation for Frontier Artificial Intelligence Models Act" (SB1047), a "modest" measure for commonsense regulation, is (as of August 2024) "the target of an extraordinary lobbying effort" "If the bill fails," he writes in the *Nation*, it will signal the "power of money in American politics" and the country's incapacity to regulate (see also Lovely 2024).

In an educational setting, chatbots also create particular challenges for student learning and academic integrity—a subject to which we now turn.

**TEACHING IDEA:** Choose one or more published resources from each of these categories for class reading and discussion, potentially by organizing the class into groups that focus on each topic. Ask the class to discuss and/or produce a set of notes on the actually existing harms of "generative AI" and the concerns such harms generate. Consider having the class draft their own recommendations for an "ethical" approach to the technology. Prepare them for the difficulty of this question perhaps by beginning with discussion of what "ethical" decision-making entails!

Consider supplementing (or succeeding) this idea by having students view the plenary panel on "Accountability and Online Safety" featuring historian Brittney Cooper and technologist Abeba Birhane at the recent DESIGN JUSTICE AI institute (begin at 29:35 on this video).

### 4. Student Learning and Academic Integrity: Research and Reflections

Despite the <u>hype over Al's supposed capacity to transform education</u>, researchers have only begun to evaluate the impact on student learning. There is, however, a century-long history of enthusiasts overpromising on the "personalized" benefits of education technology (Watters <u>2023</u>). As the *New York Times*'s Natasha Singer (<u>2025</u>)

-

<sup>&</sup>lt;sup>25</sup> See Whittaker (2021: 51), co-founder of the AI Now Institute, for the case that AI technology "cedes inordinate power" to a handful of corporations while significantly "capturing" academic research in the field. Estrin, who is the former CTO of Cisco, argues that the "hubris and determination of tech leaders to control society is threatening our individual, societal, and business autonomy." See Hao (2023) for discussion of a Stanford "transparency index" (Bommasani 2023) which, while itself arguably insufficient, found a wide range of gaps in disclosure including lists specifying the "authors, artists, and others" whose works were used for training; the use of copyrighted works; and documentation of a model's known biases and confabulations.

explains, OpenAl's current "campaign" to expand subscriptions at institutions of higher learning "is part of an escalating A.l. arms race among tech giants to win over universities and students"; moreover, Google and Microsoft "have for years pushed to get their computers and software into schools, and court students as future customers." Hence, the tech industry's current "push to A.l.-ify college education, amounts to a national experiment on millions of students." The pressure is often even greater in K-12 education: according to Alex Molnar (qtd. in Grose 2025), who directs the National Education Policy Center at the University of Colorado, generative AI is being "forced upon" K-12 schools "without any particular context or funding that would allow them to make informed decisions about what may or may not be valuable to them." Columnist Jessica Grose (2024) cites a Pew Research finding that only 6 percent of public school teachers in the US think that AI tools produce more benefit than harm.

Instead of solid evidence for Al's educational benefits, enthusiasts often assert that adoption of bots is required to prepare students for jobs. Some argue that Al will level the playing field (by equipping all students with newfound capacities), while others believe that students trained to use Al will be more employable than their peers. Few focus on teaching students *how* these technologies work—though such knowledge is integral to (critical) Al literacies. Indeed, it is by no means clear that a student trained to depend on bots for a wide range of tasks is a more attractive employee than one who has learned how to probe these tools and understand their serious limitations (e.g., Estrada 2025; Ramoni 2025a Ramoni 2025b). Certainly all students need to recognize the risks of entrusting high-stakes tasks to probabilistic statistical models: consider CNBC's May 2024 report that young job-seekers are sending companies "hundreds of the exact same cover letters word for word."<sup>26</sup>

All enthusiasts may counter that students trained to write good prompts will avoid such pitfalls. But what if the ability to write efficacious prompts—much like writing well

-

<sup>&</sup>lt;sup>26</sup> On the deleterious effects of generative AI on writing quality, see computer scientist Margaret Mitchell's January 2025 thread on Bluesky concerning platforms that are aggressively promoting the (unprompted) use of the technology for users' everyday writing tasks: these include reduction of originality, pressure to homogenize, increase of erroneous content, reduction of information diversity, deterioration of web content (and future training data), and irresponsible abuse of corporate power. See also Du, Gross, and Hong (2025) for a compelling case for prioritizing a writing process that prioritizes voice over shallow focus on surface polish [and New Yorker writer Kyle Chayka (2025) on recent studies that document the homogenizing effects of chatbots on writing in conjunction with reduced brain activity Marit MacArthur, Halm, Hall. Marit MacArthur - look for her article on the PAIRR project; but also tag her call for contributions to the Critical AI special series: https://criticalai.org/2025/07/01/cai-special-series-cfp-generative-ai-and-teaching-writing-in-higher-ed/Matthew Halm - forthcoming essay on Prompt "egineering" in Criticla AI 3.2 - coming out any day now Hall is a forthcoming essay on surveillance;]

in the first place—is a high-level sociocognitive capacity that requires experience, relevant knowledge, and control over the writing process? Many AI enthusiasts recommend that students prompt chatbots to generate a first draft, which they afterwards edit and revise. Once again, the approach may short-circuit a complex process. Whereas editing and revision require hard-won habits of critical reflection and rhetorical skill, chatbot writing suffers from simplistic, derivative, or inaccurate content. Where, then, is the evidence that students invited to skip over the tried-and-true building blocks of college writing will thrive as "prompt engineers," fact-checkers, and editors of mediocre text? Where too is the opportunity for such students to explore their own ideas and develop their voice?

In the meantime, the internet is already teeming with tips on prompts, which are warehoused on <u>websites</u>, recycled on social media, <u>surveilled by AI developers</u>, and generated <u>automatically</u> by systems including <u>ChatGPT</u> (see also Halm forthcoming). In such a milieu, students striving to learn need more than prompting techniques to develop academic work they can proudly claim as their own. Early research suggests that even "brainstorming" with bots may reduce students' confidence and "self-efficacy."<sup>2728</sup> After all, work that demonstrates creativity, thoughtfulness, care, and

\_

<sup>&</sup>lt;sup>27</sup> See Nataliya Kosmyna et al. for a study (2025) comparing LLM-users and non-users which found that the former displayed "consistent homogeneity" across a range of tasks; fell behind on ability to quote from essays composed "just minutes prior" (see also Chayka [2025] on the same study). The study concluded by predicting that LLM use is "likely" to coincide with a "decrease in learning skills." In a study of the impact of LLM-use on creativity, Harsh Kumar and colleagues (2024) observed that "participants who had no prior exposure to LLMs consistently performed better," for example by "generat[ing] more original ideas on average" than those exposed to LLMs. Their "findings suggest that while LLMs may provide short-term boosts in creativity during assisted tasks, they might inadvertently hinder independent creative performance." According to Hamsa Bastani et al. (2024) high schoolers who encountered math instruction by chatbot experienced substantial learning loss. Sabrina Habib and colleagues (2024), found that use of ChatGPT for brainstorming resulted in "reduced self-efficacy" for those still developing diverse "thinking skills" and "creative confidence": "as some participants expressed difficulty in coming up with ideas beyond what the AI offered." The study suggests that educators whose lessons enlist students to criticize, edit, and fact-check chatbot outputs may be overestimating (and thus undermining the development of) their students' competencies for these tasks. From a student perspective, ChatGPT's output of grammatically and syntactically correct prose, and its authoritative tone may seem like an unreachable ideal-not a machine-generated draft amenable to improvement from a novice writer. On the matter of coding assistants in industry, see, for example, Aslan's (2024) insider account of "the learning curve paradox." Although junior developers using coding assistants experienced some boosted productivity they "exhibited a shallow understanding of fundamental concepts. When asked why specific patterns were used, many struggled to explain their reasoning. The reliance on AI seemed to shift focus from learning to completion." (emphasis added).

<sup>&</sup>lt;sup>28</sup> Kosmyna et al. [2025] is among several recent studies to correlate chatbot use with increased homogeneity and "cognitive debt." As the Rutgers English department "Statement on Al" (n.d.) notes student learning goals in the humanities are often "carefully crafted to emphasize skills in critical thinking, research, textual analysis, and the use of evidence. That is particularly true of writing courses

resilience is hardly to be gotten at the touch of a button. While there may be real value in learning to recognize the flaws in a chatbot's outputs, that does not mean that tasking students to "write" by improving auto-generated content is a good way to inspire them, help them to cultivate their own articulacy, or sharpen their ability to think for themselves. As educators are called on to undertake experiments on their own students that the majority of business are not ready to trial on their customers or clients, the goal of preparing students for the future should not be abandoned to industry recommendations or harried administrators who have not yet had the time to develop their own critical Al literacies (Goodlad 2025).

In what follows we offer a **list of five critical AI literacies** designed for students but also useful for educators and all citizens:

- understanding the difference between AI as depicted in science fiction and actually existing technologies that might be called AI;
- understanding what "generative AI" is (including how the large language models [LLMs] on which that technology is built work, how the datasets for "training" them are collected, why these are "probabilistic" systems, what kinds of human reinforcement they require; what harms and risks to people and the environment they entail);

\_

that aim to develop habits of reading and writing that students need to meet rhetorical challenges creatively and to take risks intellectually. Learning goals in literature courses emphasize the ability to evaluate and critically assess sources and use the conventions of attribution and citation correctly, as well as to analyze and synthesize information and ideas from multiple sources to generate new insights. Depending on how they are used, generative AI tools can undermine all these goals." The statement goes on to describe the teaching of critical AI literacies as a process of "equipping students with the necessary knowledge for exercising judgment about whether or how to use these imperfect and, so far, largely untested commercial technologies" including the understanding of how such "tools work, what they are capable of, and how to contend with their ethical implications." For a set of potential learning objectives for writing courses that focus on resistance, see McIntyre, Fernandes, and Sano-Francini (n.d.). As writing center director Jane Rosenzweig writes in a helpful blog post: "It is crucial that we teach our students to think critically about generative Al-and asking them to engage with AI tools in different contexts across different disciplines will be an important part of that process. But rather than simply asking students to turn to the chatbot for 'feedback' or for any other step of the writing process, we should be helping our students understand how LLMs are trained, what types of data they are trained on, what we don't know about that data, and how bias is baked into these systems." Audrey Watters points out (2025) that disengagement among students is worsened by standardized testing and associated curricular changes; social media and other screens; the decline of reading (including the decline of parents reading aloud); the lingering effects of the pandemic; the instrumentalist narrative of education as a vehicle for "job skills"; increasing costs; growing economic inequality and growing competition for good jobs in a dwindling economy.

- learning to distinguish AI hype from useful information about new technologies
- distinguishing text generating chatbots from conventional "search engines" and understanding the various problems of using chatbots for research
- understanding the dangers of the ELIZA effect, especially in conjunction with systems that are implemented as conversational chatbots that use first person pronouns (and other anthropomorphizing language) and are marketed as companions, mentors, tutors, therapists, etc.

#### • The Economic Objectives of Generative AI

Teachers of critical AI literacies invariably encounter the tensions between the goals of higher education, and those surrounding the design and implementation of generative AI. Whereas education aims to strengthen students' articulacy, understanding, and application of knowledge, promoters of generative AI aim believe that the same commercial technologies through investors hope to accrue profit and market share, will also transform the economy–perhaps even usher in a <u>fourth industrial revolution</u>. College writing, which is a core proficiency for undergraduate education, is at the very center of this tension. Long considered a recursive *process*, college writing usually begins with reading and/or research; <u>proceeds to "pre-writing" practices such as "freewriting" and "brainstorming</u>"; and culminates in revision. By contrast, for AI developers writing is a *product*, the speedy delivery of which can maximize productivity. "The most important thing that technological advancement does," writes <u>one MIT researcher</u> on ChatGPT, is to enable workers to "produce economic output more efficiently."

Ironically, generative AI is stumbling because it has yet to deliver anything like such world-historical efficiencies—a point Goldman Sachs emphasized in a June 2024 report titled "Generative AI: Too Much Spend, Too Little Benefit?." As journalist Mateo Wong writes in a cogent analysis, the "industry is asking the world to engage in something like a trillion-dollar tautology." That is, "AI's world-transformative potential justifies spending any amount of resources, because its evangelists will spend any amount to make AI transform the world." According to reporting from September 2025, OpenAI continues to operate a tremendous loss, projecting burning \$115 billion through 2029 because of the high costs of compute and new models. Bain Capital's September 2025 technology report forecasted that meeting the "instatiable" computer

power and electricity demands of generative AI "could require \$500 billion in annual spending on new data centers. This "building rush," according to the <u>Wall Street Journal</u>, "is effectively a mega-speculative bet that the technology will rapidly improve, transform the economy and start producing steady profits." With no evidence that gen AI is improving worker productivity or businesses' return on investment, *The Economist* has begun to report on an <u>AI Trough of Disillusionment</u>."<sup>29</sup>

This situation has made the education market ever more important to tech investors. That is, if generative AI can be portrayed as a productivity boost for teachers and a trustworthy "copilot" for students, schools could provide the industry with the legitimacy, growth, and profitability that investors crave (e.g., Singer [2025], Goodlad [2025], Crano [forthcoming]).

While there is no single way to instill critical AI literacies, educators committed to such teaching, we urge, should steer clear of the industry's self-interested technodeterminism and irresponsible hype. Too often educators both in K-12 and higher ed are asked to assume that since investors are spending billions to make commercial tools accessible to students—all while embedding these tools into devices and platforms in ways that make them difficult to avoid—their primary role is now to teach the "ethical" and "responsible" use. This flawed thinking enlists teachers to whitewash the cognitive, social, and environmental impacts of an underregulated technology, while simultaneously striving to ensure student learning in the face of tools that were not designed with education in mind. The truth is that no teacher or student can neutralize Al's pervasive harms—which spring from a concentrated political economy bent on expanding profitable forms of resource-intensive surveillance, data-capture, and platform dominance. A critical AI literacies approach responds by helping students to "get the facts," and equipping them to make decisions about whether or how to use chatbots from positions of knowledge, citizenship, and

<sup>&</sup>lt;sup>29</sup>See also Sheryl Estrada's August 2025 *Fortune* article, "MIT report: 95% of generative Al pilots at companies are failing," describing a survey finding that only 5% of such pilots were accelerating revenue; Grant Grose's August 2025 *CIO* article, "GenAl Descends into Disillusionment;" Sri Ripidi's September 2025 *Information* article "OpenAl Says Its Business Will Burn \$115 billion Through 2025," which reports increased losses of more than \$80 billion since the first projection; and Bryan McMahon's September 2025 article in *The American Prospect*, reflecting on the widely perceived failures of OpenAl's release of GPT-5 in August 2025, warns that the "financial bet" on so-called artificial general intelligence (AGI) is so big (and so dependent on OpenAl's hype that failure could cause an economic depression.

<sup>30</sup> On cognitive harms see also above, including note 26. Discussion of "responsible" and/or "ethical" use in such discourse typically focuses on transparency of usage (such as the citation of chatbot content to document its use). Those who advocate for active adoption and use of chatbots in their pedagogy while simultaneously emphasizing the importance of teaching the harms of these systems, in effect reduce the question of "ethical" use to a potential requirement to acknowledge harms in the hope that they will improve over time. Critical Al literacies as we understand it, must meet a higher bar.

care. One approach to this effort emphasizes how a student's own research into how generative AI works—for example, <u>through simple "probes"</u> and/or <u>audits</u> of models—can provide skills and insights into the technology's strengths and limitations (of the kind that potential employers may value) without turning students into habitual users.

As generative AI struggles to find a firm foothold in business, professional work, and everyday life, those teaching critical AI literacies will need to distinguish hype from reality. To be sure, machine learning programs trained on high-quality data in dialogue with experts and community stakeholders can produce valuable tools in specific domains—welcome technologies that will likely be called "AI." But "generative AI"—especially when operationalized in the form of anthropomorphized chatbots—is a very different technology that, despite its high costs, persistent unreliability, and manifold harms, strives to be all things to all people all of the time.

That is why we urge resisting the pressure to succumb to technodeterminist and tautological thinking. As the suggestions for syllabi and assignments which follow make clear, teaching critical AI literacies does not entail "banning" AI, policing students, or fueling panics of any kind. Nor is the point to shame students who have been encouraged to use generative AI (especially given that such encouragement increasingly comes not only from ed tech and its online influencers and low-quality media but also from official policies adopted at some campuses). Rather, teaching critical AI literacies involves enabling students to understand what generative AI—or any other automated technology—can and cannot do. It emphasizes a student's need for informed straight-shooting that can prepare young learners to exercise judgment, and counter the hard sell and hype. Educators already know how to do the rest. 32

**TEACHING IDEA**: Choose one or more learning goals from your syllabus and invite students to discuss their ideas about how best to achieve these objectives. If these goal(s) bear directly on generative AI, invite your students to explore potential impacts on their learning. Consider inviting your students to co-create a contract or set of policies that involve possible use of permissible digital technologies. For example, should students adopt a preferred search engine or relevant library application for their

\_

<sup>&</sup>lt;sup>31</sup> Less auspicious than well-scoped, special-purpose technologies in, say, weather prediction or drug applications are machine learning systems devoted to harmful predictive technologies including algorithmic pricing, facial recognition, and (supposed) fraud detection.

<sup>&</sup>lt;sup>32</sup>See *Critical AI*'s two-part special issue, "Beyond Chatbot-K: Large Language Models, Generative AI, and the Rise of Chatbots." We have also begun an ongoing series, "Generative AI and Teaching Writing in Higher Ed,", guest-edited by Marit MacArthur (UC, Davis), the first entries of which are forthcoming in *Critical AI* 3.2 (October 2025).

research? If they are permitted to use grammar check should they commit to disabling any "generative" or "AI" features that the application offers? Ask students to specify how they believe these choices will affect learning goals and the classroom community.<sup>33</sup>

# 5. Suggestions for Updating Assignments and Syllabi (through Clarifying Course Policies and Learning Goals)

Despite an already full workload, <u>every instructor should review their syllabus and</u> <u>assignments to assure the clearest possible articulation of policies regarding the use</u> <u>of generative AI</u>. This process could begin with a close look at your institution's code of conduct and the learning goals for your course; you might also wish to <u>distribute</u> <u>our student guide to your class</u> (possibly discussing it as a classroom activity).

The code of conduct at Rutgers states that students must ensure "that all work submitted in a course, academic research, or other activity is the student's own and created without the aid of impermissible technologies, materials, or collaborations." In creating policies on generative Al tools, this puts special emphasis on the identification of permissible technologies and the question of whether a given tool impedes the learning goals of the course (including the submission of suitable work that is "the student's own").

At Rutgers, <u>learning goals vary widely across and within schools, disciplines, majors, pedagogical approaches, and levels of difficulty</u>. This means that course policies and teaching approaches that effectively build critical AI literacies may look different.

#### For example,

A computer science instructor teaching an introductory course may wish to
prohibit students from using chatbots for coding in order to ensure that they learn
fundamental skills; but she may wish to allow such use in an advanced course
designed for those who have already mastered these skills.

<sup>33</sup> No doubt there is a great deal of research to be done on how different campuses are rolling out different tools, with or without meaningful shared governance, through what campus entity, and through what kind of messaging to faculty, students, or both.

- A professor might organize "probing" experiments that enable students to investigate model bias, perhaps preparing them to publish their results on a class website.<sup>34</sup>
- An instructor teaching research at the graduate level may wish to allow students
  to use chatbots to improve grammar and syntax so long as they document that
  the actual research is their own. To save energy and avoid inadvertent changes to
  content, she may instruct her students to disable the generative features on a tool
  such as Grammarly.
- A humanities instructor who assigns writing and research to build critical thinking and sharpen engagement with course materials and themes may determine that use of chatbots poses a serious impediment to these objectives. He may therefore explain why the use of these tools is impermissible for assigned writing. However, he may simultaneously assign a research task in which students compare and contrast resources they find using search engines, library databases, or Wikipedia to those they find through chatbot use. He may also recommend that students use <a href="kagi.com">kagi.com</a> or <a href="DuckDuckGo">DuckDuckGo</a>, search engines that, unlike Google, do not surveil users or monetize their data. (Note that as of 2025, both have begun to incorporate some AI features, thus limiting the extent to these search engines help users to avoid the recourse to resource-intensive computations for simple queries.
- An environmental science course may focus on the <u>use of water</u>, energy, and rare-earth metals required to train and deploy chatbot systems without using these tools in class.
- A course in law or in graphic design may involve the study of "AI" copyright
  infringement across different companies and domains while inviting students to
  use these tools to audit models for research regarding intellectual property.

Of course, students will have their own views on the topic:

 Some may wish to opt out of using tools known to embed harmful stereotypes and/or subject users to surveillance and data collection.

<sup>&</sup>lt;sup>34</sup> For examples of simple "probing techniques" that teachers can introduce into their class, see Daniel Estrada's slides, "<u>Teach Your Students to Audit.</u>" and Teresa Ramoni's slides, "<u>Fostering Critical Al Literacy Through Research and Probing.</u>"

- Some may seek to hone their ability to use "Al" in order to prepare themselves for the job market. They may find that a professor who assigns the probing of model bias and inaccuracies has prepared them to demonstrate their skills more deeply than a course that simply allows chatbot use for the generating of text.
- Some may recognize the virtue of "probing" models but have privacy concerns.
   They may request an alternative assignment that does not require them to sign up for a surveillant tool.

The good news is that all of these situations can effectively teach and enhance critical Al literacies, whether Al tools are directly used or not.

Below we offer recommendation on ASSIGNMENTS, USE OF AI "DETECTORS," PEER REVIEW PRACTICES, and SYLLABUS UPDATES:

ASSIGNMENTS: In reviewing assignments, instructors may wish to implement changes in light of the fact that students may be tempted to use AI tools even if they are told not to do so. Simple response papers ("what did you think of this reading?") might work best in a classroom setting in handwritten fashion (or with wifi disabled on computers, phones, and tablets for those needing accommodations).

- <u>In place of conventional take-home essays</u>, consider alternatives that may boost student interest and enhance opportunities for active and project-based learning (e.g., in-class presentation, video, podcast, e-portfolio, oral examination)
- For <u>research papers and other time-intensive take-home writing</u>, try the following recommendations:
  - a) invite students to develop a topic they care about so as to encourage intrinsic motivation;
  - b) develop rubrics that emphasize critical thinking, problem-solving, applied knowledge, and use of evidence and which de-emphasize the summarization of content and the perfection of grammar and syntax. Include requirements for <u>use and citation of specific evidence, whether drawn from course materials or from independently researched sources.</u>
    (Bot-generated text tends toward summary and generality with little or no quotation; when prompted to provide quotations, bots often deliver quoted material that is fabricated or incorrect). Consider expanding the

- criteria by requiring explicit engagement of class themes or discussion (so as to avoid pat and predigested content).
- c) consider assigning a "high-stakes" pre-writing assignment (that counts for a significant part of the student's grade); this can take the form of a detailed "questionnaire" that calls on students to address specific components of the writing process that apply to course materials directly and are not readily automated.<sup>35</sup>
- d) consider demoting the importance of grammar and syntax for take-home writing in your grading rubrics. Accentuate the importance of engaging core questions, meeting specified criteria, developing ideas, using and appropriately citing concrete evidence, cultivating persuasiveness and individual voice.
- e) consider the use of a mandatory student certification comparable to that which follows: "In concert with [Rutgers'] code of conduct, which mandates "that all work submitted in a course, academic research, or other activity is the student's own and created without the aid of impermissible technologies, materials, or collaborations," this course has been designed to promote your learning, critical thinking, skills, and intellectual development without reliance on unauthorized technology including chatbots and other forms of "artificial intelligence" (AI). [Although you may use search engines, spell-check, and simple grammar-check in completing your assignments, ] you will be asked to submit your written work with the following statement. "I certify that this assignment represents my own work. I have not used any unauthorized or unacknowledged assistance or sources in completing it including free or commercial systems or services offered on the internet, or text generating systems embedded into software." Please consult with your instructor if you have any questions about the permissible use of technology in this class.
- Consider combining take-home assignments (or replacing them) with in-class exams (written by hand or with wifi disabled for students with accommodations).
- Bear in mind that student success with all kinds of assignments often depends on <u>careful reading of assigned materials and/or the taking of notes</u>. Adobe pdf readers in the effort to upsell AI features to their users now include banners like

<sup>&</sup>lt;sup>35</sup> For more details on this suggestion, see point #6 in these slides, "Teaching and Generative AI."

that pasted in below which invite readers to "save time" by reading an Al-generated summary (ignoring the document flaws of Al summaries [e.g., see DeRewal 2025] and implying the equivalence of reading any summary to reading "a long document."

This appears to be a long document. Save time by reading a summary. Generative AI User Guidelines

To help ensure careful reading, consider assigning (interesting) multiple choice exams and/or in class close reading exercises that help to encourage rigorous engagement of key passages and challenging ideas. NOTE: Many instructors find that students are surprisingly unfamiliar with multiple choice, which, if properly conceived, can help students with comparison and contrast and other important analytical skills (e.g., which of these answers is false?).

Many AI tools profess to reduce the "cognitive load" of note-taking by automating the generation of notes. Apart from the intrusion on classroom privacy that such systems present, the logic behind this technology ignores that the "cognitive load" of note-taking is part of why the practice helps students to actively learn. Consider assigning each day's note-taking (by hand or with a wi-fi- disabled device) as an assignment; although each student will take their own notes daily, the day's assigned note-taker will share their notes with the entire class.

**TEACHING IDEA**: In advance, generate an automated summary of an assigned reading (or find one online). Before showing the summary, ask students to choose one or more passages in the reading that they found especially challenging, interesting, provocative, or illuminating. After discussion, show them the summary and invite them to discuss what the automated summary identified and what it passed over, decontextualized, or misinterpreted (bear in mind that because language modeling works through pattern-finding in large datasets, summaries are often generated by locating the parts of the assigned reading that are similar to texts observed in the training data). Next (or as a separate plan) compare the synthetic summary to a relevant Wikipedia page on the text, author, or topic in question. Invite your students to compare the information available on a Wikipedia page, and the protocols for its provenance (including crowd-sourcing, citation, moderation) If—as often happen— the autogenerated summary resembles the Wikipedia entry (perhaps even repeating some of its language), discuss the ethical ramifications of this lack of credit and attribution.

Ask students how they would feel if their work on a blog, article, or wikipedia page showed up in a generative Al output without any attribution and reduced traffic to their page accordingly. Consider sharing some passages from this article which argues that Al's use of content for training data has dissolved the web's "social contract."

Al "DETECTION": The media discourse around student "cheating" is permeated by hype. Moreover, such discourse sometimes portrays college assignments as if they were task-specific labor disconnected from learning and the application of critical thinking. We believe that intrinsic motivation is one of the best ways to ensure a student's engagement with written work, research, and other forms of assessment.<sup>36</sup> A focus on cheating or plagiarism, on the other hand, can undermine the relationship between teachers and students. That said, we recognize that instructors need to ensure fairness and academic integrity in the classroom and that non-permitted use of generative AI has created added hurdles.

Nonetheless, please be aware <u>no current system being marketed to "detect" machine-generated text is reliable: false positives and false negatives are possible and even likely</u>. Some of these tools evince <u>biases against non-native English speakers</u>.<sup>37</sup> Finally, use of AI detection software, which is not FERPA-protected, may also violate students' privacy or intellectual property rights.<sup>38</sup> We suggest that instructors avoid these systems or at least discount them as reliable evidence for violations of academic integrity.

Instructors who <u>suspect the unauthorized use of an AI tool in their course</u> should consider asking for a meeting to discuss the student's approach to completing the work and refer cautiously to the problematic content. An instructor who simply asks the student to describe the writing process that led to the work in question may be more successful than one who explicitly brings up potential misuse of AI. If impermissible use of AI seems likely (as when fabricated quotations and other

<sup>36</sup> See Lang (2013) for research that examines academic dishonesty and factors that "encourage" students to engage in this behavior. Lang also has three part series in the *Chronicle of Higher Education* on this research. (Part 1, Part 2, Part 3).

<sup>&</sup>lt;sup>37</sup> See Liang et al. (2023) for specific details about this study. Anecdotal accounts of outputs circulating on social media suggest that neurodivergent people may also be at risk for discriminatory assessment. See Verma et al. (2023) for news of a new and allegedly more reliable detector created by Berkeley NLP researchers (2023).

<sup>&</sup>lt;sup>38</sup> In this context, it is worth noting that many of the digital tools students use voluntarily or according to instructor guidelines involve breaches of privacy and IP rights, including, at various points, Grammarly, Google Docs, and Zoom: on this changing landscape see, for example, Knowles (2023) and Merchant (2023).

inexplicable content shows up in a take-home assignment), a good first step might be to consult with an appropriate administrator for advice to learn more about recommended policies and suggested next steps.

If you continue to suspect misuse of generative tools, consider experimenting with changes in assignments as described above. For further conversation, feel free to reach out to <a href="mailto:criticalai@sas.rutgers.edu">criticalai@sas.rutgers.edu</a>

**PEER REVIEW:** Peer review is a worthwhile practice that helps to create community and idea-sharing among students. Unfortunately, some students may submit their classmates' work to AI systems in order to simulate peer review. This is a serious violation of academic integrity since it involves another student's intellectual work. If peer-review is an assigned take-home task, consider prefacing the assignment with a notice about the gravity of this offense and/or have students sign a statement certifying that they have not used unauthorized devices (see above for an adaptable template).

**SYLLABUS**: Whether an instructor wishes to build in the use of chatbots for certain assignments, allow students to experiment with them as they wish, or prohibit their use, we recommend clarifying these policies on syllabi and discussing them with students. Explain how you reached a decision that comports with the learning goals for the course. Consider discussing how chatbots work and the various problems described on this webpage, in concert with our <u>Student Guide</u> (see additional resources below), Colleagues at Critical AI @ Rutgers are available to answer specific questions or suggestions for teaching.

#### I. For instructors who do not want students to use Al tools for their course

When specifying on one's syllabus that the use of chatbots and other AI tools is *not* permissible, instructors should be as clear as possible and may wish to refer to the Rutgers code of conduct, cited above, in doing so. Given that AI tools are now widely incorporated seamlessly into platforms including Google, Adobe, grammar-checking tools such as Grammarly, and software suites such as Microsoft Office, a clear and specific statement is the best possible way to communicate with your students. In addition, you may wish to ask students to submit a statement of academic integrity along with their assignments.

For example,

In concert with Rutgers' code of conduct, which mandates "that all work submitted in a course, academic research, or other activity is the student's own and created without the aid of impermissible technologies, materials, or collaborations." this course has been designed to promote your learning, critical thinking, skills, and intellectual development without reliance on unauthorized technology including chatbots and other forms of "artificial intelligence" (AI). [Although you may use search engines, spell-check, and simple grammar-check in crafting your assignments,] you will be asked to submit your written work with the following statement. "I certify that this assignment represents my own work. I have not used any unauthorized or unacknowledged assistance or sources in completing it, including free or commercial systems or services offered on the internet or text generating systems embedded into software." Please consult with your instructor if you have any questions about the permissible use of technology in this class.

Below is some alternative or additional language for syllabi which was developed at the <u>University of Toronto</u>.

- The use of generative AI tools or apps for assignments in this course, including tools like ChatGPT and other AI writing or coding assistants, is prohibited.
- The use of generative artificial intelligence tools and apps is strictly prohibited in all course assignments unless explicitly stated otherwise by the instructor in this course. This includes ChatGPT and other AI writing and coding assistants.
   Use of generative AI in this course may be considered use of an unauthorized aid, which is a form of cheating.

# II. For instructors who wish to permit use of Al tools <u>in particular</u> circumstances

When specifying on one's syllabus that the use of chatbots and other AI tools is permissible in certain circumstances, instructors should be as clear as possible and may wish to refer to the Rutgers code of conduct, cited above, in doing so. Bear in mind that students may be using these tools for different purposes in different classes so that it is important to be specific in describing the particular usages you allow or encourage. Given that AI tools are now incorporated seamlessly into platforms such as Google Docs, grammar-checking tools such as Grammarly, and software suites such as Microsoft Office, a clear and specific statement that lays out permissible usages is the best possible way to communicate with your students.

For example, an instructor who does not want AI tools to be used in conjunction with written work but who wants to encourage students to do probing research on model content might consider the following statement:

In concert with Rutgers' code of conduct, which mandates "that all work submitted in a course, academic research, or other activity is the student's own and created without the aid of impermissible technologies, materials, or collaborations," this course has been designed to promote your learning, critical thinking, skills, and intellectual development without reliance on unauthorized technology including chatbots and other forms of "artificial intelligence" (AI). [Although you may use search engines, spell-check, and simple grammar-check in crafting your assignments,] you will be asked to submit your written work with the following statement. "I certify that this assignment represents my own work. I have not used any unauthorized or unacknowledged assistance or sources in completing, it including free or commercial systems or services offered on the internet or text generating systems embedded into software."

A partial exception to this policy is an authorized [exploration of model bias which we will conduct in Week X in order to build your learning on critical AI literacies.]

Please consult with your instructor if you have any questions about the permissible use of technology in this class.

(As above, our recommendation is that any instructor assigning work that involves mandatory use of an Al tool consider developing an option for students who have data privacy or other concerns.)

### III. For instructors who wish to permit use of AI tools

When specifying on one's syllabus that the use of chatbots and other AI tools *is permissible* (*or assigned*), instructors should be as clear as possible about how this decision comports with the learning goals for their course and may wish to refer to the Rutgers code of conduct, cited above, in doing so. Instructors may also want to emphasize critical AI literacies including the importance of recognizing that current AI tools are subject to bias, misinformation, environmental harms *et al.* (as discussed above). Given the widespread availability of a variety of tools, be sure to be clear and specific about which tools are permitted and, if applicable. what forms of citation are required to document such use.

For example,

In concert with Rutgers' code of conduct, which mandates "that all work submitted in a course, academic research, or other activity is the student's own and created without the aid of impermissible technologies, materials, or collaborations," this course has been designed to help you develop knowledge about the use and abuse of Al tools. Al tools may be used as an aid in the creative process, but with the understanding that this should be accompanied by independent evaluation, critical thinking, and reflection. Students who choose to use these tools are responsible for any errors or omissions resulting from their use. They will also be required to provide as an appendix the prompts used, the generated output, and a thoughtful reflection on the outcomes. When appropriate, students may also be asked to consider the environmental and social costs of using the tool.

(As above, our recommendation is that any instructor assigning work that involves mandatory use of an AI tool consider developing an option for students who have data privacy or other concerns.)

Some instructors who permit use of AI tools for written assignments implement syllabus statements like these, developed at the <u>University of Toronto</u>.

- Students must submit, as an appendix with their assignments, any content produced by an artificial intelligence tool, and the prompt used to generate the content.
- Students may choose to use generative AI tools as they work through the
  assignments in this course; this use must be documented in an appendix for
  each assignment. The documentation should include what tool(s) were used,
  how they were used, and how the results from the AI were incorporated into the
  submitted work.

### 6. Select Resources You Might Wish to Read or to Share with Your Students

This document already includes many resources that you might enjoy or share with students and colleagues. Here we provide some additional resources. As this is a living document, we plan to continue to update it with additional resources as they become available. Please feel free to suggest them to us. [WE ARE UPDATING THE BELOW LIST TO INCLUDE BRIEF SUMMARIES AND TO GROUP ACCORDING TO TOPIC - COMING SOON!]

- 404 Media, The 404 Media Year in Review (podcast), 31 December 2024.
- AI Now Institute, AI Now Salon Series (videos).
- Abebe, Rediet and Maximilian Kasy, "The Means of Prediction." The Boston Review, 20 May 2021.
- Abraham, Roshan. "A Black Hole of Energy Use': Meta's Massive AI Data Center is Stressing Out a Louisiana Community." 404 Media. 23 June 2025.
- Ahmad, Meher, Jessica Grose, and Tressie McMillan Cottom. "What A.I. Really Means for Learning." New York

  Times, 12 August 2025.
- Ahmed, Kalim. "The Internet is a Place Where No One Has an Accent." AI Policy Perspectives, 15 July 2025.
- Allison, Leslie and Tiffany DeRewal, "Where Knowledge Begins? Generative AI as Search and the Problem of Friction." forthcoming in *Critical AI* 2.2 (October 2024).
- Al-Sibai, Noor. "OpenAI Admits That Its New Model Still Hallucinates More Than a Third of the Time." Futurism, 1 March, 2024.
- Anaya. "AI Image Generators Often Give Racist and Sexist Results: Can They Be Fixed?" Nature, 19 March 2024.
- Ars Staff. "As Internet Enshittification Marches On, Here Are Some Of The Worst Offenders." Ars Technica. February 5, 2025.
- Aslan, Murat. "The Hidden Costs of AI Coding Assistants: Insights From a Senior Developer." Medium. 19 December 2024.
- Atleson, Michael. "Succor Borne Every Minute." Federal Trade Commission (weblog), 11 June 2024.
- "Attorney General Alan Wilson Leads 44 AGs Demanding Big Tech End Predatory AI Targeting of Kids." 2025. South Carolina Attorney General's Office. August 25, 2025.
- Baak, Stefan. "The Human Decisions that Shape Generative AI. Who is Accountable for What?" *Mozilla Foundation Newsletter*, 2 August, 2023.
- Bastian, Matthias. "Apple AI Researchers Question OpenAI's Claims About o1's Reasoning Capabilities." The Decoder. 12 October 2024.
- BBC. "Data Centre Power Use 'To Surge Six-Fold in 10 Years." 26 March 2024.
- Bender, Emily M. "ChatGPT Has No Place in the Classroom." Mystery AI Hype Theater 3000 Newsletter, 2024.

Bender, Emily M. "<u>ChatGP-why: When, if ever, is synthetic text safe, appropriate, and desirable?</u>" (video-recorded lecture for GRAILE AI), 12 August, 2023.

- Bender, Emily M., Timnit Gebru, Angelina Mcmillan-Major, and Shmargaret Shmitchell. "On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? "Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency, Virtual Event, Canada, March 2021, 610-623.
- Bender, Emily M. and Alex Hanna, "AI Causes Real Harm. Let's Focus on That vs. the End-of-Humanity Hype." Scientific American, 12 August, 2023.
- —--. Mystery AI Hype Theater 3000. (podcasts and videos).
- Bianchi, Federico, Pratyusha Kalluri, Esin Durmus, Faisal Ladhak, Myra Cheng, Debora Nozza, Tatsunori Hashimoto et al. "Easily Accessible Text-to-Image Generation Amplifies Demographic Stereotypes at Large Scale." *ArXiv* preprint 2022.
- Birhane, Abeba, Atoosa Kasirzadeh, David Leslie, and Sandra Wachter. 2023. "Science in the Age of Large Language Models." *Nature Reviews Physics* 5, 277–80. April 26, 2023.
- Birhane, Abeba and Deborah Raji, "<u>ChatGPT, Galactica and the Progress Trap.</u>" *Wired*, 9 December 2022.
- Bjork, Collin. "Clones in the Classroom: Why Universities Must Be Wary of Embracing AI-Driven Teaching Tools." *The Conversation*. 18 September 2024.
- Blackwell, Alan. "Ooops! We Automated Bullshit (ChatGPT is a Bullshit Generator. To Understand AI, We Should Think Harder about Bullshit)" 11 September, 2023.
- Blunt, Katherine. "Who Pays? AI Boom Sparks Fight Over Soaring Power Costs." Wall Street Journal, 29 July 2025
- Bode, Katherine. "Why You Can't Model Away Bias." MLQ 81. no 1 (2020) 95-124.
- Booth, Robert. "Paul McCartney Warns AI 'Could Take Over' As UK Debates Copyright Laws." *The Guardian*. 10 December 2024.
- Broussard, Meredith. <u>Artificial Unintelligence: How Computers Misunderstand the World.</u> Cambridge: MIT Press, 2019.
- Burke, Garance, and Hilke Schellmann. "Researchers Say an AI-Powered Transcription Tool Used in Hospitals

  Invents Things No One Ever Said." AP News. 26 October 2024.
- Challapally, Aditya, Chris Pease, Ramesh Raskar, and Pradyumna Chari. "<u>The GenAI Divide: State of AI In Business 2025</u>." *MIT Project Nanda*. July 2025.
- Conrad, Kathryn. "A Blueprint for An AI Bill of Rights for Education." Critical AI (blog), n.d. [Sneak Preview for an upcoming April 2024 issue of Critical AI]
- Costanza-Chock, Sasha. <u>Design Justice: Community-Led Practices to Build the Worlds We Need</u>. Cambridge: MIT Press, 2020.
- Cotton, Tressie McMillan. "The Tech Fantasy That Powers A.I. Is Running on Fumes." New York Times. 29 March 2025.

Crabapple, Molly. "Molly Crabapple – Debates in AI – RISD." Lecture presented at Rhode Island School of Design's Debates in AI Symposium, Providence, RI, April 2024.

#### Data Workers

Inquirihttps://docs.google.com/document/d/14Fp1DbmjZx9lzz8kvUo57yDZ9I6l5v4quLX6tbIpPvQ/edit?tab=t.0es: Invited Pieces.

Critical AI (@CriticalAI). "Yes @OpenAI new text-to-video is impressive but here's 5 questions that journos & the public should be asking–relentlessly. 1. When will you release the". X, February 16, 2024, 9:22am. <a href="https://twitter.com/CriticalAI/status/1758497093172207811">https://twitter.com/CriticalAI/status/1758497093172207811</a>

da Silva, João. "Google Turns to Nuclear to Power AI Data Centres." BBC. 15 October 2024.

DeRewal, Tiffany. "Ask an Expert: Evaluating LLM 'Research Assistants' and Their Risks For Novice Researchers." Critical AI (blog), n.d.

Doctorow, Cory. "The Enshittification of TikTok." Wired. 23 January 2023.
——. "Enshittification' is Coming for Absolutely Everything." Financial Times 8 February 2024.

Drhal, Carmen. "AI Was Asked To Create Images of Black African Docs Treating White Kids. How'd It Go?."

NPR. 6. October, 2023.

Economic Security Project. "New Shadow Report Addresses Gaps in the Senate Roadmap for AI and Calls for a Policy Agenda That Harnesses AI in the Public Interest." 5 May 2024. Economic Security Project Report. Accessed June 25, 2024.

Ellis, Lindsay, and Katherine Bindley. "AI Is Wrecking an Already Fragile Job Market for College Graduates." Wall Street Journal, 28 July 2025.

Estrada, Daniel. "Teaching Insights: How to Teach AI to Students: AI Ethics and NJIT Audit Project." Critical AI Blog. (n.d.)

Estrin, Judy. "The Case Against AI Everything, Everywhere, All at Once." Time, 11 August 2023.

Farrell, Maria and Robin Berjon. "We Need to Rewild the Internet." Noema. Berggruen Institute, 16 April 2024.

Flaherty, Colleen. "How AI Is Changing—Not 'Killing'—College." Inside Higher Ed, 29 August 2025.

Flynn, Shannon. 2020. <u>"The Difference between Symbolic AI and Connectionist AL."</u> RE•WORK Blog - AI & Deep Learning News. September 24, 2020.

Gilliard, Chris. "How Ed Tech Is Exploiting Students." Chronicle of Higher Education, 8 April 2018.

Goldman Sachs. "AI is poised to drive 160% Increase in Data Center Power Demand." May 14, 2024.

Goodlad, Lauren M. E. "Editor's Introduction: Humanities in the Loop." Critical AI (2023) 1-2.

- Goodlad, Lauren M.E. and Sam Baker, "Now the Humanities Can Disrupt AL." Public Books, 20 February, 2023.
- Goodlad, Lauren M. E., and Kathryn Conrad. March 2024. "Teaching Critical AI Literacies." NORRAG Policy Insights #4 AI and Digital Inequities,
- Goodlad, Lauren M. E., and Matthew Stone. "Beyond Chatbot-K: On Large Language Models.

  'Generative AI,' and the Rise of Chatbots; An Introduction" Critical AI (2024) 2.1
- Gross, Grant. "Devs Gaining Little (If Anything) From AI Coding Assistants." CIO. 26 September 2024.
- Guo, Eileen. "A Major Al Training Data Set Contains Millions of Examples of Personal Data." MIT Technology Review, 22 July 2025.
- Guo, Eileen. "An AI Chatbot Told a User How To Kill Himself—But the Company Doesn't Want To 'Censor' It." MIT Technology Review. February 6, 2025.
- Hadgu, Asmelash Teka, and Timnit Gebru. "Replacing Federal Workers with Chatbots Would Be a Dystopian Nightmare." Scientific American. Springer Nature. April 14, 2025.
- Hao, Karen. "Microsoft's Hypocrisy on AI." The Atlantic. 13 September 2024.
- Hao, Karen (@\_KarenHao). 2024. "For years I've been interviewing data annotation workers who are the lifeblood of the AI industry [...]." X thread, March 16, starting at 10:24 AM. <a href="https://twitter.com/">https://twitter.com/</a> KarenHao/status/1769006784273101074
- Hao, Karen. "AI is Taking Water From the Desert." The Atlantic. 1 March 2024.
- Heaven, Will Douglas. "The Algorithm." MIT Technology Review Newsletter, 25 August 2025.
- Herrman, John. "ChatGPT and Google Gemini Are Both Doomed." Intelligencer, 1 March, 2024.
- Hsu, Tiffany and Stuart A. Thompson, "<u>Disinformation Researchers Raise Alarms About A.I Chatbots.</u>" *New York Times.* 8 February 2023. Updated 20 June 2023.
- Huckins, Grace. "Why the AI Moratorium's Defeat May Signal a New Political Era." MIT Technology Review, 9 July 2025.
- Jaźwińska, Klaudia and Aisvarya Chandrasekar. "AI Search Has A Citation Problem." *Columbia Journalism Review*. Tow Center for Digital Journalism. 6 March 2025.
- Kaltheuner, Frederike, Leevi Saari, Amba Kak, and Sarah Myers West. "I. Reorienting European AI and Innovation Policy." *AI Now Institute*. 15 October 2024.
- Kak, Amba, Susan Meyers West and Meredith Whittaker. "Make No Mistake–AI Is Owned by Big Tech." *MIT Technology Review.* 5 December 2023.
- Kirmer, Stephanie. "The Cultural Backlash Against Generative AI." Medium. 1 February 2025.

- Krietzberg, Ian. "OpenAI Accuses New York Times of Paying Someone to Hack ChatGPT." *TheStreet*. 27 February 2024.
- Kumar, Harsh, Jonathan Vincentius, Ewan Jordan, and Ashton Anderson. "<u>Human Creativity in the Age of LLMs:</u>
  <u>Randomized Experiments on Divergent and Convergent Thinking.</u>" *arXiv*, 2410.03703, September 24, 2024.
- Lopatto, Elizabeth. "Stop Using Generative AI As a Search Engine." The Verge, 5 December 2024.
- Luccioni, (Alexandra) Sasha, Bruna Tevelin, and Margaret Mitchell. "<u>The Environmental Impacts of AI Primer</u>." *Hugging Face*, 3 September 2024.
- LuccioniAlexandra Sasha, Yacine Jernite, and Emma Strubell. "Power Hungry Processing: Watts Driving the The Cost of AI Deployment?" *ArXiv.* May 2024.
- Maiberg, Emanuel, and Jason Koehler. "<u>Inside the Booming 'AI Pimping' Industry</u>." *404 Media*. 20 November 2024.
- Matthewson, Tara García. "AI Chatbots Can Cushion the High School Counselor Shortage But Are They Bad For Students." *The Markup*. 4 March 2025.
- Mauran, Cecily. "The Era of the AI-Generated Internet is Already Here." Mashable. 27 January, 2024.
- McKendrick, Joe. "Generative AI May Be Creating More Work than It Saves." ZDNet. 24 May, 2024.
- Merchant, Brian. "The AI Industry Has a Battle-Tested Plan to Keep on Using Our Content without Paying for It."

  The LA Times. 12 January 2024.
- Metz, Cade, Cecilia Kang, Sheera Frenkel, Stuart A. Thompson, and Nico Grant. "How Tech Giants Cut Corners to Harvest Data for A.I." New York Times. 6 April 2024.
- Mills, Anna, Maha Bali, and Lance Eaton. "How do we respond to generative AI in education? Open educational practices give us a framework for an ongoing process." 6. 1. (2023).
- Monserrate, Steven Gonzalez. "The Staggering Ecological Impacts of Computation and the Cloud." MIT Press Reader. 14 February 2022.
- Muldowny, Decca, and Alex Hanna. "Sora 2 Serves Up More Slop." Mystery AI Hype Theater 3000: The Newsletter. 10 October, 2025.
- Narayanan, Arvind and Sayash Kapoor, AI Snake Oil (substack)
- Nicoletti, Leonardo and Dina Bass. "<u>Humans Are Biased. Generative AI Is Even Worse.</u>" *Bloomberg Technology.* 12 June 2023.
- Niemeyer, Kenneth. "Billionaire Larry Ellison Says a Vast AI-Fueled Surveillance System Can Ensure 'Citizens Will Be On Their Best Behavior." Business Insider. 15 September 2024.

- O'Brien, Isabel. "<u>Data Center Emissions Probably 662% Higher Than Big Tech Claims. Can It Keep Up the Ruse?</u>" *The Guardian.* 15 September 2024.
- O'Donnell, James. "AI's Giants Want to Take Over the Classroom." MIT Technology Review, 15 July 2025.
- O'Neil, Lorina, "These Women Tried to Warn Us about AI." Rolling Stone, 12 August 2023.
- Pasquale, Frank. New Laws of Robotics: Defending Human Expertise in the Age of AI. Cambridge: Harvard University Press, 2020.
- Peters, Uwe, and Benjamin Chin-Yee. "Generalization Bias in Large Language Model Summarization of Scientific Research." Royal Society Open Science, vol. 12, no. 4, 2025.
- Petri, Alexandra. "<u>I Hate the Gemini 'Dear Sydney' Ad More Every Passing Moment</u>." *The Washington Post.* 31 July 2024.
- Pierce, David. "The text file that runs the internet" [on how training AI has broken the "social contract" of the Internet.] *The Verge.* 14 February 2024.
- Ramkumar, Amrith and Brian Schwartz. "Silicon Valley Launches Pro-AI PACs to Defend Industry in Midterm Elections." *Wall Street Journal*, 25 August 2025.
- Reiley, Laura. "What My Daughter Told ChatGPT Before She Took Her Life." New York Times, 24 August 2025.
- Rosalsky, Greg. "10 Reasons Why AI May be Overrated." NPR, 6 August 2024.
- Rosenzweig, Jane. "<u>Teaching Insights: What Happens When a Novice Writer Asks ChatGPT for Editing Advice?</u>" *Critical AI Blog* (n.d.)
- Sanders, Nathan, and Bruce Schneier. "AI Is Changing How Politics Is Practiced in America." *The American Prospect*, 10 October 2025.
- Saul, Josh, Naureen S. Malik, and Mark Chediak. "<u>AI Boom Is Driving a Surprise Resurgence of US Gas-Fired Power</u>." *Bloomberg*. 16 September 2024.
- Schmidt, Eric, and Selina Xu. "Silicon Valley is Drifting Out of Touch with the Rest of America." New York Times, 19 August 2025.
- Shumailov, Ilia, Zakhar Sumaylov, Yiren Zhao, Yarin Gal, Nicholas Papernot, Ross Anderson. "The Curse of Recursion: Training on Generated Data Makes Models Forget." *ArXiv* (self-published preprint uploaded on 31 May 2023)
- Smith, Gary and Jeffrey Funk. "When It Comes to Critical Thinking, AI Flunks the Test." The Chronicle of Higher Education. 12 March 2024.
- Stone, Matthew, "Large Language Models: a Whirlwind Tutorial." (video skip to 15:31)

  Critical AI (blog), n.d. [Recorded as part of a 2021 virtual workshop for "The

  Ethics of Data Curation," an NEH-supported research collaboration between Rutgers and

  Australian National University] (See also Stone's talk on the first panel of the videos for
  the Critical AI Literacy in a Time of Chatbots symposium.

Tacheva, Jasmina, and Srividya Ramasubramanian. "AI Empire: Unraveling the interlocking systems of oppression in generative AI's global order." *Big Data & Society* 10.2 (2023): 20539517231219241.

Tan, Rebecca and Regine Cabato. "Behind the AI boom, an army of overseas workers in 'digital Sweatshops." The Washington Post, 28 August 2023

Tiku, Natasha, Kevin Schaul, and Szu Yu Chen. "<u>These Fake Images Reveal How AI Amplifies Our Worst Stereotypes.</u>" *The Washington Post.* 1 November 2023.

Tech Won't Save Us (@techwontsaveus). 2024. "<u>DATA VAMPIRES</u> is a four-part series exploring the costs of hyperscale data centers and why [...]." X thread, October 7, starting at 10:17 AM.

Urquieta, Claudia and Daniela Dib. "U.S Tech Giants are Building Dozens of Data Centers in Chile. Locals Are Fighting Back." *Rest of World.* May 31, 2024.

Vassel, Faye-Marie, Evan Shieh, Cassidy R. Sugimoto, and Thema Monroe-White. 2024. "The Psychosocial Impacts of Generative AI Harms." *Proceedings of the AAAI Symposium Series* 3 (1): 440–47.

Vincent, James. "How Much Electricity Does AI Consume?" The Verge, February 16, 2024.

Warner, John. "Everyone Should Read 'Teaching Machines." Inside Higher Ed. 28 September 2024.

Watters, Audrey. 2025. "AI Slop Education." Second Breakfast

Watters, Audrey. <u>Teaching Machines: The History of Personalized Learning</u>. Cambridge: MIT Press, 2021.

Wiggers, Kyle. "It Sure Looks Like OpenAI Trained Sora on Game Content – And Legal Experts Say That Could Be A Problem." *TechCrunch*. 11 December 2024.

Wilkins, Joe. "AI Therapist Goes Haywire, Urges User to Go on Killing Spree." Futurism, 25 July 2025.

Williams, Rua M. (@FractalEcho). 2024. "The racism behind chatGPT that we aren't talking about... This year, I learned that students use [...]." X thread, February 17, starting at 4:37 PM. https://twitter.com/FractalEcho/status/1758968904674836979

Wong, Matteo. "The Schools Without ChatGPT Plagiarism." The Atlantic. 25 October 2024.

Wong, Matteo. "The AI Boom Has an Expiration Date." The Atlantic. 17 October 2024.

Wong, Matteo, "The Generative-AI Revolution May Be a Bubble." Atlantic Monthly, 2 August 2024.

Wong, Matteo, "America Already Has An AI Underclass." Atlantic Monthly, 26 July 2023.

Zakarin, Jordan, and Josh Hirschfeld-Kroen. 2023. "This Is the Biggest Threat to Actors & Writers." YouTube Video. *More Perfect Union*.