

---

# Balance Human Agency & AI Assistance in the Tussle for the “Right” to \*

---

Anonymous Authors<sup>1</sup>

## Abstract

AI increasingly mediates and augments daily life, yet its technical properties risk reshaping the way users choose (what they consume), own (what they create), learn, and work. These risks reflect widely recognized AI principles, which unfortunately remain largely high-level and unoperationalizable. As the risks arising from AI assistance continue to subtly erode human agency, this work takes a different approach to operationalizing AI principles by translating the risks into concrete research questions that can guide the design of AI systems. This enables the research community to balance human agency and AI assistance in AI development and align it with the goal of benefiting humans.

## 1. Introduction

Artificial intelligence (AI) has become inseparably intertwined with the rhythms of modern life, shaping how we *choose* what we consume, *own* what we create, *learn* new information, and *work* to complete tasks<sup>1</sup>. In leisure, it subtly chooses the next video or reel we watch (Dekker et al., 2025). In creative production, it composes text and images that both rival and enrich traditional forms of art (Khan, 2024a). In workplaces, it lightens the burden of labor through automation (Simon, 2025), and streamlines the intricacies of office work (Leopold, 2025). When we seek understanding, it serves as an ever-present and patient tutor (Zhai et al., 2024; Kosmyna et al., 2025). Simultaneously, scholarly engagement with AI has shifted from niche to mainstream, resulting in an exponential growth of AI-engaged publications across

<sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

<sup>1</sup>This paper focuses on *choosing*, *owning*, *working*, and *learning* because they are universal and pervasive in everyday life. Furthermore, the consequences of failing to safeguard these capabilities or “rights” are profound and demand timely and coordinated action. We discuss the selection of these capabilities or “rights” in Section A.3.

diverse scientific fields and applications (Duede et al., 2024) to support the development of capable assistants that benefit mankind (Li, 2025) by augmenting our abilities to choose, own, work, and learn.

However, there are also social ramifications to AI use as certain technical properties of AI systems can reshape the way its users choose, own, work, and learn. These technical properties reflect widely recognized principles at the core of frameworks and guidelines on AI use (Jobin et al., 2019; Floridi & Cowls, 2022) regarding human autonomy, dignity, fairness, and meaningful participation<sup>2</sup>. While it is clear that AI research should engage with these principles, existing work highlights a gap between these high-level principles and their practical operationalization (Birkstedt et al., 2023; Morley et al., 2021; 2023). Indeed, conceptual papers far outnumber empirical, technological or implementational research (Birkstedt et al., 2023). Hence, it remains difficult to translate these principles into concrete research directions.

This gap has direct implications on how AI impacts human agency (in choosing, owning, working, and learning), challenging not merely expectations or established routines, but the fundamental role of humans as being able to exercise meaningful discretion in areas of importance to us. Assistive technologies like AI do not merely relieve us of certain drudgeries, they risk fostering overreliance beyond cognitive offloading (Zhai et al., 2024) and skill decay (Noy & Zhang, 2023). What is at stake is more than our capacity to act—it is in fact our ability to make decisions. Crucially, human agency is the capacity to act, but it exists only when it is being exercised (Sanchini et al., 2019). AI risks undermining human agency not merely through automation or the intrinsic opacity of black-box algorithms, but through system outputs shaped by influences that are simultaneously technical and institutional, and largely inaccessible to users. As intrinsic opacity, expert control, and organizational mediation become embedded in these outputs, the scope for exercising individual agency is progressively reduced (Winner, 1980).

In line with AI research’s mission to build systems that benefit and augment humans, this paper takes a different approach to operationalizing these AI principles. It examines

<sup>2</sup>We discuss mapping these widely recognized AI principles to the “rights” in a framework in Section A.2.

how AI creates risks to everyday activities such as choosing, owning, working, and learning and translates these into concrete research questions that can guide the design and evaluation of AI systems, enabling the research community to better align AI development with its mission. We adopt the vernacular of “rights”<sup>3</sup> to identify the capabilities we risk giving up to our machine assistants: our “rights” to choose, to own, to work, and to learn.

To start the ball rolling, we present open research questions to address to reinforce each “right”. These research questions distill reflections from researchers<sup>4</sup> actively engaged in AI development. The questions are deliberately output-oriented and grounded in lived practice, and act as technical metrics specifying the intended behavior/capability of an AI system, but presented in plain language for a wider audience who may consider addressing these questions from the perspective of their field<sup>5</sup>. Each question is accompanied by an appendix that surveys relevant prior work, highlights unresolved technical challenges that render the question non-trivial, and presents illustrative technical case studies or solutions that serve as inspirations (rather than prescriptions) with the aim of motivating (potentially better) alternative solutions. Throughout the paper, we use “AI systems” as a broad term encompassing, but not limited to, generative models/AI assistants, AI agents, computer vision, voice or speech AI assistants, decision-making systems, recommendation systems, AI productivity tools, personalization systems, and other tools built on AI. We introduce specificity in system types only when necessary to preserve question generality, particularly when discussing risks. We hope these questions will help the audience strike a balance between human agency and AI assistance so that society as a whole reaps the transformative rewards of AI.<sup>6</sup>

## 2. The “Right” to Choose

**RISK: Echo Chambers, Filter Bubbles, and Information Barriers.** AI models can increasingly align and optimize their responses with human expectations (Christiano et al., 2017; Rafailov et al., 2023). Large pretrained generative

<sup>3</sup>The “rights” framing is not intended to propose new legal human rights, but to align with established AI ethics principles synthesized in global surveys on normative expectations of users (Jobin et al., 2019). We conceptualize “rights” in Section A.1.

<sup>4</sup>In August 2025, a cohort of [redacted]-funded Master’s and PhD scholars convened to deliberate on the risks associated with AI use. As researchers engaged in developing AI systems, they support AI development and seek to maximize its transformative potential for broad societal benefit.

<sup>5</sup>This target audience includes researchers, developers, social scientists, humanities scholars, and policymakers, who can work together to co-design AI systems to benefit society.

<sup>6</sup>Amid debates over the proper scope of AI agency (Bengio et al., 2025; LeCun, 2022), this paper examines how to preserve human agency without framing it as a zero-sum tradeoff with AI.

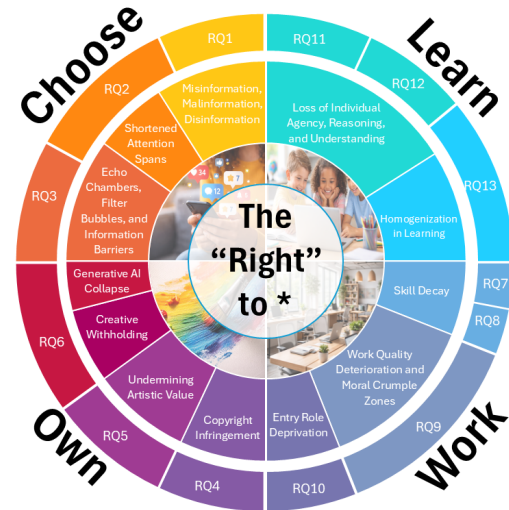


Figure 1. Overview: We focus on the risks (inner ring) arising from AI adoption (without careful scrutiny) that threaten the four “rights” to choose, own, work, and learn. We propose corresponding research questions (RQs in outer ring) to mitigate these risks. Each RQ is color-mapped to one risk, but some span multiple risks and hence extend across risk boundaries (e.g., RQ3). Note: Slice sizes are illustrative and do not indicate risk or RQ importance.

models curate suggestions based on user preferences, search histories, or their own latent biases (Liang et al., 2021), while recommendation feeds curate the content (Tik Tok, 2025) instead of relying on explicit user selection (Agan et al., 2023; Mills & Sætra, 2024). This narrows exposure to diverse perspectives, reinforcing confirmation biases and echo chambers and trapping individuals in filter bubbles (Anwar et al., 2024). Over time, users may struggle to articulate their information needs (Hirvonen et al., 2024), and blocking their own access to broader knowledge.

**RISK: Shortened Attention Spans, Cognitive Offloading and Cognitive Decline.** AI-driven platforms maximize engagement through short, personalized, and rapidly delivered content. As users scroll through AI-curated streams of brief videos, attention spans become increasingly fragmented and decisions become increasingly simple. For example, YouTube presents users with a list of videos to choose from (Ghori et al., 2025), while TikTok further simplifies this decision into binary choice by presenting one video at a time (Tik Tok, 2025). Over time, repeated exposure conditions users’ attention systems to favor immediate, frictionless AI curation over diversity or depth of choices (Shanmugasundaram & Tamilarasu, 2023). Individuals may struggle with activities that require sustained cognitive focus, such as long-form reading, and turn to brainrot content (Yazgan, 2025; Yousef et al., 2025), reflecting a decline in their cognitive abilities.

This phenomenon is observed beyond AI-driven platforms. The use of LLMs over traditional search engines for informa-

tion gathering has also been found to ease cognitive load but limits the depth and quality of content engagement (Stadler et al., 2024). This reveals the cost of cognitive ease. Rather than freeing users to think more deeply, LLMs instead reduce the need to engage with the information independently and meaningfully.

**RISK: Mis/Mal/Disinformation.** AI systems enable large-scale generation and amplification of fabricated media and coordinated false narratives, undermining public discourse integrity<sup>7</sup> (Maitland & Lee, 2025). Such campaigns “firehose” or flood information ecosystems with high-volume content to overwhelm users (Paul & Matthews, 2016). Empirical work further underscores how users face difficulty in discerning AI-generated content (Frank et al., 2024; Groh et al., 2024), undermining the capacity of users to make informed choices about what to believe.

The risks mentioned above are non-exhaustive and will evolve with the way we access information. AI systems risk obscuring the “Right” to Choose as they increasingly mediate everyday decisions. Safeguarding this “right” therefore requires preserving informed agency throughout the decision process. For example, by ensuring users have the freedom to access or reject unbiased or desired sources of information or options, and even retaining the ability to review, override, or withdraw consent for choices delegated to AI systems.

**RQ1. How can AI systems present diverse and credible content that expands users’ informational exposure, while preserving meaningful and sustained engagement?**

AI systems act as choice architects, structuring and presenting limited sets of options<sup>8</sup> for users to choose from (Agan et al., 2023; Mills & Sætra, 2024). However, optimizing option sets solely for engagement metrics such as click- or watch-time risks placing users in echo chambers (Baumann et al., 2024) and increases their vulnerability to misinformation (Avram et al., 2020). This has motivated regulations to protect users to have their personal data such as cookies, age, or geographical location *forgotten* (European Parliament and Council of the European Union, 2016) from AI systems, reduce excessive personalization, and allow users to experience more general and unbiased content. In addition to legal enforcement, more can be done in terms of AI development. Consider, for example, platforms that adopt multi-objective recommendation frameworks that explicitly trade off engagement with viewpoint diversity/credibility and civic value. Domain experts and journalists can help

<sup>7</sup>As part of the influence operation: Storm-1516, AI-generated artwork and synthetic personas were used to stage false claims. These narratives were laundered through a network of AI-enabled news websites and algorithmically boosted by automated accounts on platforms like Telegram and X, demonstrating how AI systems can support coordinated, large-scale disinformation campaigns.

<sup>8</sup>We formulate an example option set in Section B.1.1.

identify informational blind spots in various fields, assess credibility, and distinguish factual content from opinion. This can support informed and autonomous choices about what users engage with and expand users’ informational exposure while supporting meaningful and sustained engagement, allowing users to escape filter bubbles and avoid undue exposure to false information. As mentioned in Section 1, every RQ, such as RQ1, specifies the intended AI system behavior, invites interdisciplinary examination, and is accompanied by an appendix describing related work, open challenges, and illustrative technical solutions; for RQ1, see Section B.1.2.

**RQ2. How can we design AI systems that users can trust?**

The opacity of AI systems prevents humans from fully understanding and interpreting its workings, even for computer scientists or those with specialized training (Burrell, 2016). Increasing transparency is a common response to opacity, but trust depends on appropriate calibration, as both under- and over-transparency can erode trust in the AI system’s recommendations, especially when users have previously disagreed with the AI system’s outputs (Kizilcec, 2016). Collaborating between legal scholars, social scientists, domain experts, along with computer scientists can help to build AI systems that justify the reason or intent behind their recommended option set. With interpretations for AI outputs, users can ensure that the intentions behind each option—at least at the level of interpretation—are sufficiently informative, and also unbiased, constructive, engaging, and free from harmful influence. This will allow users to trust and evaluate options critically and make autonomous and well-reasoned decisions on their own terms, without worry of manipulation or incorrectness. See Section B.1.3.

**RQ3. How can AI systems empower users to reject, adjust, or withdraw from algorithmic curation?**

Even when the option set is optimal and the decision is informed, users may still wish to “fight the feed” (Ghori et al., 2025) by adjusting or withdrawing consent to algorithmic curation, especially as their preferences evolve. AI platform developers may support this by providing the appropriate user controls. Alternatively, AI systems can also be developed to directly support users in wiping out their preferences from the platform.

Algorithmic intervention is complex and its preferred extent varies between users. Even determining the “correct” amount of algorithmic intervention is inherently difficult as users may accept personalized services but simultaneously reject the collection of personal data required for it (Kozyreva et al., 2021). The acceptance of algorithmic curation also depends on social norms, emotions, and heuristics (Acquisti et al., 2015). These challenges have motivated regulations on consent withdrawal (European Parliament

and Council of the European Union, 2016) and works on preference shifts (Wang et al., 2023); further progress requires collaboration with decision scientists in psychology, economics, and sociology to design AI systems that account for shifting preferences and perspectives, and empower user control over personalization and privacy needs that evolve over time. Collaboration with public policy and legal experts is also necessary to design AI systems that re-balance power asymmetry between users (who may be naive, uncertain, and vulnerable) and governments and corporations that design and operationalize algorithmic curation. See Section B.1.4.

### 3. The “Right” to Own

Generative AI has demonstrated remarkable capability to compose text (Brown et al., 2020; Sutskever et al., 2011), images (Goodfellow et al., 2014; Rombach et al., 2022), and audio (Engel et al., 2017; van den Oord et al., 2016), and has been praised as democratizing art by turning passive observers into creators as they no longer need years of training to recreate Rembrandt’s brushstrokes or Beethoven’s virtuosity (Park, 2025). Text-, image-, audio-based works, and their mixed-media combinations are traditionally covered by copyright where a single discernible author owns each work (U.S. Copyright Office). This contrasts with the layered process of *AI-mediated* (i.e., AI-generated, AI-augmented, or created through human-AI collaboration) creation, involving artists/authors/creators, model developers, data curators, and prompters, causing AI-generated art to fall into a grey area of copyright infringement.

**Definition 3.1. Ownership** refers to the possession of a creation and the control over the reproduction, distribution, and derivative use of the creation.

Section 3.1 begins with the premise that a rose by any other name would smell as sweet: AI-mediated creative work can still be art. We focus on what follows, namely the risks to ownership or attribution of art, brought on by recognizing AI-mediated creative work as art, and focus on questions that alleviate these risks. Section 3.2 considers how when attribution fails or is not done, the “gardeners” of art (i.e., artists or authors) are not acknowledged and may lose motivation to cultivate novel art. This necessitates questions that protect “gardeners” so that we can continue to appreciate new blooms.

#### 3.1. A Rose by Any Other Name Would Smell as Sweet

**RISK: Copyright Infringement.** Recognizing AI-generated art as forms of creative expression<sup>9</sup> exposes it to

<sup>9</sup>It is unproductive to draw hard lines around what counts as “art”, especially as artists themselves increasingly use AI tools to create. We draw parallels to Benjamin’s analysis of mechanical

a wide range of infringement risks. Generative models may make unauthorized copies of copyrighted works for training or reproduce them as unlicensed derivatives, thereby violating reproduction rights and moral rights. They can also embed copyrighted materials within model weights, leading to secondary infringement when models are shared (Chesterman, 2024; Jamali, 2025; Montgomery, 2025). These practices undermine creators’ rights and recognition.

**RISK: Undermining Artistic Value.** By automating the reproduction of creative expression without attribution, AI systems appropriate human-created art, reducing demand for human artists and displacing opportunities for salaried employment in creative industries (Chesterman, 2024; Senftleben, 2023; UN Conference on Trade and Development, 2024). As generative AI tools receive recognition, this undermines human artistic integrity and the sustainability of creative labor, threatening the professional recognition and employment opportunities of artists.

#### RQ4. How can the authors of an AI-mediated artwork be identified and their contributions quantified to assign credit and ownership?

Intellectual property is a legal framework designed to incentivize and reward only human creativity and innovation (Chesterman, 2024). AI-generated or modified art falls into a legal gray area, and human and AI modifications may be interleaved throughout the creative process, making it difficult to recognize and reward authors, much less assign ownership. In reasoning about ownership of an AI-mediated creation, attribution and credit become a necessary conceptual precursor as ownership depends on understanding who has contributed what to the final work. To do so requires first identifying contributors, then quantifying their contributions, and lastly assigning credit and ownership.

One may consider the following solution. First, identifying all contributors to a piece of art acknowledges those who may claim authorship. This provides a foundation for AI researchers to trace its complete chain of authorship and influence, and is necessary as artists may use AI tools for assistance or modification, full generation, or not at all. Secondly, the magnitude and nature of each contribution can be quantified through collaboration with artists. Third, discussion on assigning credit welcomes insights from legal experts and economists, such as game theorists, who can consider each contribution and build frameworks to assign credit and ownership accordingly. This rewards creativity and innovation and assigns ownership, thus upholding artistic integrity, sustaining creative labor, and allowing continued creative contribution. See Section B.2.1.

reproduction (Benjamin, 2018) and shift attention from definitional debates to how changes in production transform a work’s aura and its real-world impact (e.g., on human perception, social relations, politics, power, misinformation, further discussed in Section 2).

220 **RQ5. How can intangible contribution be recognized**  
 221 **and attributed in AI-mediated artworks?**

222 Style is the technique of creating artwork, while expression  
 223 encompasses creative choices that convey intent, emotion,  
 224 or narrative. Even minor pixel-level changes can subtly shift  
 225 symbolism or emotional tone, altering intangible meaning  
 226 and perceived value. Legally, expression is copyrightable,  
 227 but style is not (WIPO, 2025; WTO, 1995). However, generative  
 228 AI allows style to be replicated accurately and at  
 229 scale, raising questions about whether style should now  
 230 be copyrighted to protect original creators. Accordingly,  
 231 RQ5.1 asks whether AI-mediated stylistic influence constitutes  
 232 expressive contribution and deserves creative credit. If  
 233 such stylistic influence is recognized, credit can be assigned  
 234 to the ones creating the AI-mediated artwork. Otherwise,  
 235 credit reverts solely to the original creator. RQ5.2 asks how  
 236 the extent of that intangible contribution can be measured  
 237 to reward the ones being attributed to it.

238 **RQ5.1. How can we determine when AI-mediated**  
 239 **stylistic influence deserves creative credit?**

240 Style is not protected under copyright law, but certain styles  
 241 are so distinctive that they are immediately associated with  
 242 a specific creator. Generative models can easily copy style  
 243 and blur the line between inspiration and replication. An  
 244 example is the “Studio Ghibli style”, which remains legally  
 245 unprotected (SCMP, 2025). Methods can be developed  
 246 alongside artists to quantify expressive vs. stylistic contribu-  
 247 tions. This help to distinguish creative contribution from super-  
 248 ficial reproduction and recognize creative influence that  
 249 shapes expression in AI-mediated art. In parallel, engaging  
 250 legal experts can help with acknowledging contributions  
 251 and determining copyright infringement. See Section B.2.2.

252 **RQ5.2. How can we quantify the expressive overlap**  
 253 **between an AI-mediated artwork and the original?**

254 While style replication by generative AI is often appar-  
 255 ent (SCMP, 2025), the replication of expression is less  
 256 straightforward. For example, the New York Times alleges  
 257 in an ongoing lawsuit that OpenAI is infringing on copy-  
 258 rights by reproducing “significantly more expressive content  
 259 from [an] original article than what would traditionally be  
 260 displayed” by an online search (Pope, 2024; NYT, 2023).  
 261 Having a measure that can quantify expressive overlap (like  
 262 continuous measures such as similarity scores in plagiarism  
 263 detection systems (Turnitin, 2025)) can help generative AI  
 264 rebuke or artists prove claims of copyright infringement.  
 265 Artists, legal experts, and AI developers can collaborate to  
 266 build legal frameworks based on this measure, and impose  
 267 graduated penalties based on the extent of infringement to  
 268 properly credit the original creator. See Section B.2.3.

3.2. A World Without Gardeners

If credit is attributed solely to the final visible creator such  
 as someone writing a prompt, the layered contributions of  
 other creators risk being overlooked. Authors who see their  
 creations used without recognition risk becoming disincenti-  
 vated to produce or share new works. Over time, erosion  
 of credit may discourage artistic participation, diminishing  
 the diversity and vitality of the creative ecosystem.

**RISK: Creative Withholding.** Art shared online is easily  
 copied, modified, or incorporated into datasets without  
 proper attribution, and routinely scraped from public plat-  
 forms by generative models. Many artists have expressed  
 dissatisfaction with the use of their work in AI training with-  
 out consent (Ali & Breazeal, 2023; Kambur & Dolunay,  
 2024; Lovato et al., 2024), becoming increasingly reluctant  
 to publicly share their art (Ali & Breazeal, 2023).

**RISK: Generative AI Collapse.** Creative withholding re-  
 sults in the lack of new training data. When models are  
 trained recursively on their own outputs, their quality (pre-  
 cision) or diversity (recall) progressively decrease (Alemo-  
 hammad et al., 2024; Briesch et al., 2024; Shumailov et al.,  
 2024). If human authors stop producing new work, genera-  
 tive models may spiral into self-consumption.

RQ6. **What can authors do to protect their works prior**  
**to legal judgment?**

Authors today rely on legal interpretation and subjective  
 judgment to prove that their works have been copied,<sup>10</sup> and  
 have no tangible technical mechanism to detect or demon-  
 strate overlap (de Leon, 2024). Authors may also have had  
 their works “scraped” (lawfully or unlawfully, with or with-  
 out permission) without being compensated for it (Chester-  
 man, 2024). Developing models that allow legal experts  
 to obtain tangible evidence of unauthorized use (whether  
 directly by copying or indirectly for model training) of art-  
 works will empower authors to privately verify and docu-  
 ment potential misuse, giving them awareness and control  
 before formal legal action, and encourage continued creation  
 instead of creative withholding. See Section B.2.4.

4. The “Right” to Work

**RISK: Skill Decay.** In contemporary workplaces, AI as-  
 sistants can automate routine tasks, boosting efficiency and  
 task completion (Dell’Acqua et al., 2023; Noy & Zhang,  
 2023). However, as users increasingly rely on AI tools to  
 reduce the information processing requirements of a task  
 so as to reduce cognitive demand, cognitive offloading re-

<sup>10</sup>Singaporean photographer Zhang Jingna sued Luxembourg  
 artist Jeff Dieschburg who had painted her 2017 Harper’s Bazaar  
 Vietnam photograph and submitted it to an art competition, win-  
 ning a prize. The case initially ruled against her and this was only  
 overturned two years later in 2024 (de Leon, 2024).

sults (Risko & Gilbert, 2016; Grinschgl et al., 2021; Cavicchi et al., 2025), causing the relevant human competencies to atrophy from sporadic use (Macnamara et al., 2024; Natali et al., 2025). For highly skilled users, AI tools even distract by giving generic suggestions (Brynjolfsson et al., 2023) that “relieve” users from practising or developing key cognitive processes. Over time, these diminish the ability of employees to contribute in the workplace without AI assistance, contributing to skill decay.

**RISK: Work Quality Deterioration and Moral Crumple Zones.** AI tools especially benefit less experienced workers (Brynjolfsson et al., 2023; Cazzaniga et al., 2024), allowing them to shift efforts towards more impactful tasks like planning, evaluation, and problem solving (Filippucci et al., 2024). However, this creates an illusion of competence, where workers over-estimate their own performance on the task (Fernandes et al., 2026). This also gives workers the confidence to submit hallucinated text (Ashktorab et al., 2025) or even “workslop”—long, fluent, AI generated text that seems polished, but in fact needs heavy correction—(Skibba, 2026), as work deliverables.

Furthermore, generic AI tools also produce generic solutions (Anderson et al., 2024; Brynjolfsson et al., 2023). Their widespread deployment hence creates a “homogenization” effect that reduces diversity and creativity in output, undermining work quality. As accountability is typically concentrated with the end user, workers may be placed in a moral crumple zone (Elish, 2019) (i.e., made to absorb the force of impact in a car or AI system’s accident) and unproductively blamed for negative outcomes arising from system behaviors or hallucinations they cannot fully control or detect.

**RISK: Entry-Role Deprivation.** AI tools now perform many basic tasks traditionally assigned to entry-level staff (Leopold, 2025), directly competing with human graduates for entry-level positions (Teng et al., 2024). As fresh graduates lose hiring opportunities to AI tools, they miss out on training and development opportunities entry-level roles traditionally provide. Over time, they lack the deep expertise and specialized knowledge needed to make sustained contributions to their workplaces (Nolan, 2025).

The risks mentioned above are non-exhaustive and will continue to evolve as AI systems increasingly perform human tasks. However, the intrinsic value of contributing meaningfully at work remains key to workers’ motivation and the production of original and high-quality output (Bailey & Madden, 2016; Martela, 2023). It is hence vital to ensure that workers can continue to invest effort and even modulate or decline AI assistance in AI mediated workplaces. We propose the following RQs.

**RQ7. Should AI systems collaborate or automate?**

Even when users have domain expertise, appropriate AI support can improve performance (Vaccaro et al., 2024). The use of AI systems in the workplace can improve productivity, but the balance between collaboration and automation is a trade-off where efficiency gains from automation often come with limitations (Brynjolfsson et al., 2023). This trade-off has to account for how labor should be divided between human and AI (Vaccaro et al., 2024). RQ7 decomposes this trade-off into sub-RQs that examine when (RQ7.1) and how (RQ7.2) AI can intervene.

**RQ7.1. How can AI systems detect/learn optimal moments to intervene in a task (vs. leaving it to the human)?**

One may consider AI systems that intervene in human tasks to collaborate. Identifying when to intervene requires insights from psychology to, for example, understand user cognition and recognize when human limits are reached. Domain expertise can also characterize when task difficulty or error accumulation crosses a learned threshold for intervention. Crucially, there are also situations where humans outperform AI, and no intervention is necessary, as discussed in Section B.3.6. Sparse and purposeful intervention supports problem solving while ensuring that humans remain primarily responsible. This preserves agency and mitigates skill decay. See Section B.3.1.

**RQ7.2. How can AI systems intervene in ways that complement their users’ efforts and expertise?**

By collaborating with psychologists to translate theories of cognition and learning into system-level design constraints, AI developers can calibrate the extent and mode of AI assistance to complement the user. For example, by allowing the AI system to form a contextual understanding of the user’s task, prior actions, and work environment, it can offer appropriate guidance. Adaptive and socially intelligent intervention strategies, coupled with constraints on the granularity of assistance, engage users in the problem-solving process. This preserves human agency, improves task performance, and ensures that users retain competence, expertise, and confidence for future tasks. See Section B.3.2.

**RQ8. What strategies can help organizations ensure that AI is used appropriately and sparingly so that employees continue to develop their own skills?**

The overuse of AI in the workplace can lead to skill decay and hinder skill development (Macnamara et al., 2024; Natali et al., 2025). So, frameworks can be created for companies to detect overuse of AI tools on tasks that do not require AI assistance or integrate “friction” into their AI adoption processes. For example, by drawing on behavioral and learning sciences to define appropriate task difficulty and skill retention criteria, companies can introduce small, well-defined, and appropriately difficult tasks for employees to complete without AI assistance. Employees found using AI can undergo targeted retraining to ensure their core skills

remain intact and prevent skill decay. See Section B.3.3.

**RQ9. How can AI systems support fair responsibility allocation in human–AI collaborative work?**

When AI systems go wrong and responsibility must be allocated, shrinking moral crumple zones requires legal and organizational frameworks that distribute accountability across developers, managers, and organizations rather than focus blame on the final human user (Elish, 2019). This calls for collaboration between legal and technical experts to formalize an AI system’s limits and uncertainty and reallocate accountability from end users to upstream actors.

Simultaneously, to reduce the occurrence of wrong outputs being used in the workplace, AI systems need to recognize and communicate uncertainty. By drawing on psychologists’ insights into human reasoning and logicians’ formal accounts of inference, AI systems can, for example, be designed to identify sources of uncertainty arising from epistemic limitations. Communicating uncertainty can also alert users to unverifiable or hallucinated outcomes. This helps users avoid moral crumple zones and think beyond the model’s response, fostering higher quality, diverse, and contextually grounded work outputs. See Section B.3.4.

**RQ10. How can AI systems enable a junior employee to reach the competency of a senior employee?**

As firms increasingly view AI systems as viable substitutes for fresh graduates, many become disincentivized from hiring junior employees (Teng et al., 2024). This restricts fresh graduates from accessing and gaining a foothold in organizations. To mitigate this entry-role deprivation, the capability development of junior employees can be accelerated so that their performance becomes comparable to senior employees. In highly specialized domains such as software engineering, urban planning, and manufacturing, companies may prefer candidates with 10-20 years of experience over those with only 1-2 years. One may consider AI systems that can patch this skill gap and endow junior employees with hard and soft skills typically accumulated through long-term practice. This helps junior employees avoid entry-role deprivation and enables them to learn on the job and approach senior-level proficiency at an earlier stage by accelerating the development of their cognitive maps. See Section B.3.5.

While Section 4 protects the “Right” to Work, it does not claim preference for human over AI effort. Section B.3.6 acknowledges scenarios where AI may outperform humans.

**5. The “Right” to Learn**

**Definition 5.1. Learning** (Gross, 2012) is the process of acquiring new understanding, knowledge, behaviors, skills, values, attitudes, and preferences.

**RISK: Teacher-Learner AI Divide.** In conventional educa-

tional settings, learners learn from and follow the guidance of teachers. AI disrupts this relationship by giving learners access to alternative sources of instruction beyond the teacher’s oversight. Students may use AI to learn different “truths” or copy homework from. Teachers are therefore forced to play catch-up with students’ AI use, making it harder for them to understand how learning actually takes place. Without insight into this learning process, teachers may struggle to engage, diagnose learning needs, and adapt their teaching, which can ultimately harm student performance (Rowan & Grootenboer, 2016).

**RISK: Reduced Confidence in Education.** In conventional educational settings, syllabuses, certified teachers, and standardized assessments provide some assurance that learners are reaching the right outcomes and engaging in a meaningful process of learning. With AI as a middleman in this process, strong standardized test performance does not necessarily guarantee that the learner has understood the material as they may offloaded reasoning to the AI or sidestep the meaningful path to understanding learning materials. This is especially so as educational technologies frame learning as “effortless” (Mentutor, 2024; Shabanov, 2025). This undermines the value of syllabuses, standardized testing and learning outcomes, and erodes confidence that people, especially learners, teachers, and pedagogists may have in education itself<sup>11</sup>.

These risks mentioned above will continue to evolve, especially as future learners leverage AI as an encyclopedic “compression” of the internet for learning and AI tools for education continue to evolve<sup>12</sup>. We introduce research questions (RQs) to mitigate the risks above and safeguard the “Right” to Learn with AI for learners who use AI responsibly<sup>13</sup> to enhance rather than bypass learning.

**RESEARCH QUESTION 11 (RQ11). How can AI systems be designed to close or even bridge the gap between teachers and learners?**

Teachers should be able to see the learner’s questions and responses from an AI system, as these reveal what learners struggle with and what instructions they have received. Teachers should also be able to control the complexity and scope of answers that an AI gives, such as by limiting the concepts that an AI can use. This would help align AI sup-

<sup>11</sup>Other problems also arise from having AI as an intermediary in learning. Risks such as the loss of the learner’s individual agency, reasoning, and understanding and homogenization in learning are presented in Section B.4.1

<sup>12</sup>Section B.4.7 discusses preliminary progress by ChatGPT and Gemini in addressing RQs in Section 5.

<sup>13</sup>Learners may begin with genuine and sincere intent but gradually drift towards AI dependence due to the ease of obtaining ready-made answers. We hope the risks discussed above will ground sincere learners in their commitment to learn meaningfully. We elaborate on this in Section A.1.

port with upcoming lessons and account for syllabuses or differences across educational contexts. For example, in learning mathematics, teachers may prefer to teach specific mental models such as algebra, model-drawing (Ng, 2022), or intuitive reasoning. By bridging the AI gap between the teacher and learner, teachers can understand students’ learning progress, empathise with, and get a greater knowledge of the individual needs of their students, to support their learning (Zhou, 2022). See Section B.4.2.

#### RQ12. How can we design AI systems that guarantee meaningful learning?

One solution places the responsibility for meaningful learning primarily on the learner. AI systems that act like teachers can be designed to measure the depth of the learner’s understanding and dynamically calibrate their support (RQ16), strategically allow learners to make mistakes (RQ17), or adapt to preferred learning strategies (RQ18). Here, the learner remains responsible for engaging with AI in ways that preserve genuine understanding instead of bypassing the learning process.

The alternative solution places the responsibility on the teacher. Here, AI systems should give teachers oversight and control over AI use, and ensure that AI supports pedagogical goals, like in RQ11.

Ultimately, meaningful learning is personal and context-dependent. The best arrangement may be a blend of solutions, where both learners and teachers are accountable for ensuring that AI supports meaningful learning. See Section B.4.4 to B.4.6.

## 6. Conclusion

AI assistants increasingly mediate daily human processes. While they may offer short-term gains, they also risk hollowing out the capacities that sustain human agency. These risks reflect widely recognized AI principles, which unfortunately remain largely high-level and unoperationalizable. We therefore emphasize the importance of operationalizing these principles through concrete research questions that translate daily human processes into design choices in AI systems. These questions are intentionally output-oriented and grounded in lived practice, but presented in plain language for wider audiences including researchers, developers, social scientists, humanities scholars, and policymakers. We hope that by championing research answering these questions, humans can retain a core set of concrete capacities and preserve agency when choosing, owning, working, and learning in an AI-mediated world.

## 7. Alternative Views

An alternative to our “rights”-based framing is the principle-based framing of AI ethics, which focuses on operationalizing high-level ethical principles such as fairness, accountability, and transparency into technical systems. Building on the convergence of such principles identified in surveys (Jobin et al., 2019), a growing body of work has sought to translate these principles into practice (Hagendorff, 2020; Raji et al., 2020).

While the principle-based approach is valuable, our work is intended as complementary. Rather than beginning from desired attributes of AI systems, our approach foregrounds domains of human experience that we argue are worth preserving. This reframing shifts the focus from system-level properties to human-centered outcomes, and provides a different form of guidance for technical development, as reflected in our output-oriented research questions.

At the same time, a “rights”-based framing opens up the potential for a further line of research: these proposed “right” can, in future work, be more explicitly anchored to broader existing human rights frameworks, whose initial framing may not necessarily be grounded in the context of AI. This creates a pathway for linking technical research questions to established legal and institutional mechanisms for rights protection, thereby situating AI development within broader governance structures.

Additionally, we address other potential questions (e.g., practical feasibility of answering RQs, exhaustiveness of the “rights” and RQs) the reader may have in the FAQ (Section A.1).

## References

- Acquisti, A., Brandimarte, L., and Loewenstein, G. Privacy and human behavior in the age of information. *Science*, 347(6221):509–514, 2015.
- Agan, A. Y., Davenport, D., Ludwig, J., and Mullainathan, S. Automating automaticity: How the context of human choice affects the extent of algorithmic bias. Working Paper 30981, National Bureau of Economic Research, 2023.
- Agarwal, D., Naaman, M., and Vashistha, A. AI suggestions homogenize writing toward western styles and diminish cultural nuances. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*, pp. 1–21, 2025.
- Agarwal, K. AI in education – evaluating ChatGPT as a virtual teaching assistant. *International Journal for Multidisciplinary Research (IJFMR)*, 5(4), 2023.
- Ahmad, S. F., Han, H., Alam, M. M., Rehmat, M. K., Irshad, M., Arraño-Muñoz, M., and Ariza-Montes, A. Impact of artificial intelligence on human loss in decision making, laziness and safety in education. *Humanities and Social Sciences Communications*, 10(1):311, 2023.
- Aldasoro, I., Gambacorta, L., Korinek, A., Shreeti, V., and Stein, M. Intelligent financial system: How AI is transforming finance. *Journal of Financial Stability*, 81: 101472, 2025.
- Alemohammad, S., Casco-Rodriguez, J., Luzi, L., Humayun, A. I., Babaei, H., LeJeune, D., Siahkoochi, A., and Baraniuk, R. Self-consuming generative models go MAD. In *Proceedings of the International Conference on Learning Representations*, 2024.
- Ali, S. and Breazeal, C. Studying artist sentiments around AI-generated artwork. arXiv:2311.13725, 2023.
- AllSides. AllSides Media Bias Chart, 2025. URL <https://www.allsides.com/media-bias/media-bias-chart>. Accessed: 2025-11-03.
- Amnesty International. Taiwan becomes first in Asia to legalize same-sex marriage after historic bill passes, 2019. URL <https://www.amnesty.org/en/latest/press-release/2019/05/taiwan-same-sex-marriage-law/>. Accessed: 2026-01-03.
- Anderson, B. R., Shah, J. H., and Kreminski, M. Homogenization effects of large language models on human creative ideation. In *Proceedings of the 16th conference on creativity & cognition*, pp. 413–425, 2024.
- Anderson, L. W. and Krathwohl, D. R. *A taxonomy for learning, teaching, and assessing: A revision of Bloom’s taxonomy of educational objectives: complete edition*. Addison Wesley Longman, Inc., 2001.
- Anwar, M. S., Schoenebeck, G., and Dhillon, P. S. Filter bubble or homogenization? Disentangling the long-term effects of recommendations on user consumption patterns. In *Proceedings of the ACM Web Conference 2024*, pp. 123–134, 2024.
- Ashktorab, Z., Desmond, M., Pan, Q., Johnson, J., Brachman, M., Dugan, C., Danilevsky, M., and Geyer, W. Emerging reliance behaviors in human-AI content grounded data generation: The role of cognitive forcing functions and hallucinations. In *Proceedings of the 4th Annual Symposium on Human-Computer Interaction for Work*, 2025.
- Avram, M., Micallef, N., Patil, S., and Menczer, F. Exposure to social engagement metrics increases vulnerability to misinformation. *Harvard Kennedy School Misinformation Review*, 2020.
- Bagchi, C., Menczer, F., Lundquist, J., Tarafdar, M., Paik, A., and Grabowicz, P. Social media algorithms can curb misinformation, but do they? arXiv:2409.18393, 2024.
- Bailey, C. and Madden, A. What makes work meaningful — or meaningless. *MIT Sloan Management Review*, 2016.
- Bailey, J. and Warner, J. AI tutors: Hype or hope for education? *Education Next*, 25(1):62–71, 2024.
- Baumann, F., Halpern, D., Procaccia, A. D., Rahwan, I., Shapira, I., and Wüthrich, M. Optimal engagement-diversity tradeoffs in social media. In *Proceedings of the ACM Web Conference 2024*, pp. 288–299, 2024.
- Bengio, Y., Cohen, M., Fornasiere, D., Ghosn, J., Greiner, P., MacDermott, M., Mindermann, S., Oberman, A., Richardson, J., Richardson, O., Rondeau, M.-A., St-Charles, P.-L., and Williams-King, D. Superintelligent agents pose catastrophic risks: Can scientist AI offer a safer path? arXiv:2502.15657, 2025.
- Benjamin, W. The work of art in the age of mechanical reproduction. In *Modern Art and Modernism: A Critical Anthology*, pp. 217–220. Routledge, 2018.
- Bertschinger, N., Rauh, J., Olbrich, E., Jost, J., and Ay, N. Quantifying unique information. *Entropy*, 16(4):2161–2183, 2014.
- Bhandari, A. and Bimo, S. Why’s everyone on TikTok now? The algorithmized self and the future of self-making on social media. *Social Media + Society*, 8(1): 20563051221086241, 2022.

- 495 Bhatt, U., Chen, V., Collins, K. M., Kamalaruban, P., Kal-  
496 lina, E., Weller, A., and Talwalkar, A. Learning person-  
497 alized decision support policies. In *Proceedings of the*  
498 *AAAI Conference on Artificial Intelligence*, volume 39,  
499 pp. 14203–14211, 2025.
- 500 Birkstedt, T., Minkkinen, M., Tandon, A., and Mäntymäki,  
501 M. AI governance: themes, knowledge gaps and future  
502 agendas. *Internet Research*, 33(7):133–167, 2023.
- 503 Black, P. and Wiliam, D. Assessment and classroom learn-  
504 ing. *Assessment in Education: principles, policy & prac-*  
505 *tice*, 5(1):7–74, 1998.
- 506 Bleher, H. and Braun, M. Diffused responsibility: attri-  
507 butions of responsibility in the use of AI-driven clinical  
508 decision support systems. *AI and Ethics*, 2(4):747–761,  
509 2022.
- 510 Boden, M. A. *The Creative Mind: Myths and Mechanisms*.  
511 Routledge, 2004.
- 512 Bomba, F. and De Angeli, A. Agency and authorship in  
513 AI art: Transformational practices for epistemic troubles.  
514 *International Journal of Human-Computer Studies*, pp.  
515 103652, 2025.
- 516 Brand, A., Allen, L., Altman, M., Hlava, M., and Scott, J.  
517 Beyond authorship: Attribution, contribution, collabora-  
518 tion, and credit. *Learned Publishing*, 28(2), 2015.
- 519 Brashier, N. M. and Schacter, D. L. Aging in an era of fake  
520 news. *Current directions in psychological science*, 29(3):  
521 316–323, 2020.
- 522 Briesch, M., Sobania, D., and Rothlauf, F. Large language  
523 models suffer from their own output: An analysis of the  
524 self-consuming training loop. arXiv:2311.16822, 2024.
- 525 Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan,  
526 J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G.,  
527 Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G.,  
528 Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu,  
529 J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M.,  
530 Gray, S., Chess, B., Clark, J., Berner, C., McCandlish,  
531 S., Radford, A., Sutskever, I., and Amodei, D. Language  
532 models are few-shot learners. In *Advances in Neural*  
533 *Information Processing Systems*, volume 33, pp. 1877–  
534 1901, 2020.
- 535 Brynjolfsson, E., Li, D., and Raymond, L. R. Generative  
536 AI at work. Working Paper 31161, National Bureau of  
537 Economic Research, 2023.
- 538 Buchholz, A., London, B., Di Benedetto, G., Lichtenberg,  
539 J. M., Stein, Y., and Joachims, T. Counterfactual ranking  
540 evaluation with flexible click models. In *Proceedings*  
541 *of the 47th International ACM SIGIR Conference on Re-*  
542 *search and Development in Information Retrieval*, pp.  
543 1200–1210, 2024.
- 544 Burrell, J. How the machine ‘thinks’: Understanding opacity  
545 in machine learning algorithms. *Big Data & Society*, 3  
546 (1):2053951715622512, 2016.
- 547 Cavicchi, S., Abubshait, A., Siri, G., Mustile, M., and Cia-  
548 rdo, F. Can humanoid robots be used as a cognitive  
549 offloading tool? *Cognitive Research: Principles and*  
550 *Implications*, 10(1):17, 2025.
- 551 Cazzaniga, M., Jaumotte, F., Li, L., Melina, G., Pantou,  
552 A. J., Pizzinelli, C., Rockall, E. J., and Tavares, M. M.  
553 Gen-AI: Artificial intelligence and the future of work.  
554 *Staff Discussion Notes*, 2024(001):1, 2024. ISSN 2617-  
555 6750.
- 556 Chan, G. and Lo, E. Is AI cheating on the rise? Few  
557 cases reported by S’pore universities, but experts warn  
558 of risks. *The Straits Times*, 2025. URL <https://www.straitstimes.com/singapore/is-ai-cheating-on-the-rise-few-cases-reported-by-spore-universities-but-experts-warn-of-risks>. Accessed: 2025-11-15.
- 559 Charusaie, M.-A., Mozannar, H., Sontag, D., and Samadi, S.  
560 Sample efficient learning of predictors that complement  
561 humans. In *Proceedings of the International Conference*  
562 *on Machine Learning*, pp. 2972–3005. PMLR, 2022.
- 563 Chen, R., Zhang, X., Luo, M., Chai, W., and Liu, Z.  
564 Pad: Personalized alignment of LLMs at decoding-time.  
565 arXiv:2410.04070, 2024.
- 566 Chen, Z., Niu, X., Foo, C.-S., and Low, B. K. H. Broaden  
567 your SCOPE! Efficient multi-turn conversation planning  
568 for LLMs with semantic space. In *Proceedings of the*  
569 *International Conference on Learning Representations*,  
570 2025.
- 571 Chesterman, S. Good models borrow, great models steal:  
572 intellectual property rights and generative AI. *Policy and*  
573 *Society*, 44(1):23–37, 2024.
- 574 Chi, M. T. H. Self-explaining: The dual processes of gener-  
575 ating inference and repairing mental models. In *Advances*  
576 *in instructional psychology: Educational design and cog-*  
577 *nitive science*, volume 5, pp. 161–238, 2000.
- 578 Christiano, P. F., Leike, J., Brown, T., Martic, M., Legg,  
579 S., and Amodei, D. Deep reinforcement learning from  
580 human preferences. In *Advances in Neural Information*  
581 *Processing Systems*, volume 30, 2017.
- 582 Coffelt, T. A., Grauman, D., and Smith, F. L. M. Employers’  
583 perspectives on workplace communication skills: The

- 550 meaning of communication skills. *Business and Profes-*  
 551 *sional Communication Quarterly*, 82(4):418–439, 2019.
- 552  
 553 Corral, M. The harm & hypocrisy of AI art,  
 554 2023. URL [https://www.corralldesign.com/](https://www.corralldesign.com/writing/ai-harm-hypocrisy)  
 555 [writing/ai-harm-hypocrisy](https://www.corralldesign.com/writing/ai-harm-hypocrisy). Accessed: 2025-  
 556 11-06.
- 557  
 558 Dalva, Y., Venkatesh, K., and Yanardag, P. FluxSpace: Dis-  
 559 entangled semantic editing in rectified flow transformers.  
 560 arXiv:2412.09611, 2024.
- 561  
 562 de Leon, E. A. A. Singaporean photographer wins appeal  
 563 in copyright case against artist from luxembourg,  
 564 2024. URL [https://asiaiplaw.com/sector/](https://asiaiplaw.com/sector/copyright/singaporean-photographer-wins-appeal-in-copyright-case-against-artist-from-luxembourg)  
 565 [copyright/singaporean-photographer-](https://asiaiplaw.com/sector/copyright/singaporean-photographer-wins-appeal-in-copyright-case-against-artist-from-luxembourg)  
 566 [wins-appeal-in-copyright-](https://asiaiplaw.com/sector/copyright/singaporean-photographer-wins-appeal-in-copyright-case-against-artist-from-luxembourg)  
 567 [case-](https://asiaiplaw.com/sector/copyright/singaporean-photographer-wins-appeal-in-copyright-case-against-artist-from-luxembourg)  
 568 [against-artist-from-luxembourg](https://asiaiplaw.com/sector/copyright/singaporean-photographer-wins-appeal-in-copyright-case-against-artist-from-luxembourg). Ac-  
 569 cessed: 2025-11-06.
- 570  
 571 Dekker, C. A., Baumgartner, S. E., and Sumter, S. R. For  
 572 you vs. for everyone: The effectiveness of algorithmic per-  
 573 sonalization in driving social media engagement. *Telem-*  
 574 *atics and Informatics*, 101:102300, 2025.
- 575  
 576 Dell’Acqua, F., McFowland, E., Mollick, E. R., Lifshitz-  
 577 Assaf, H., Kellogg, K., Rajendran, S., Krayner, L., Can-  
 578 delon, F., and Lakhani, K. R. Navigating the jagged  
 579 technological frontier: Field experimental evidence of the  
 580 effects of AI on knowledge worker productivity and qual-  
 581 ity. *Harvard Business School Technology & Operations*  
 582 *Mgt. Unit Working Paper*, (24-013), 2023.
- 583  
 584 Doshi, A. R. and Hauser, O. P. Generative AI enhances  
 585 individual creativity but reduces the collective diversity  
 586 of novel content. *Science Advances*, 10(28):eadn5290,  
 587 2024.
- 588  
 589 Duede, E., Dolan, W., Bauer, A., Foster, I., and Lakhani, K.  
 590 Oil & water? diffusion of AI within and across scientific  
 591 fields. arXiv:2405.15828, 2024.
- 592  
 593 Dutta, A., Mehrab, K. S., Sawhney, M., Neog, A., Khurana,  
 594 M., Fatemi, S., Pradhan, A., Maruf, M., Lourentzou, I.,  
 595 Daw, A., and Karpatne, A. Open world scene graph gen-  
 596 eration using vision language models. arXiv:2506.08189,  
 597 2025.
- 598  
 599 Elish, M. C. Moral crumple zones: Cautionary tales in  
 600 human-robot interaction. *Engaging Science, Technology,*  
 601 *and Society*, 5:40–60, 2019.
- 602  
 603 Engel, J., Resnick, C., Roberts, A., Dieleman, S., Norouzi,  
 604 M., Eck, D., and Simonyan, K. Neural audio synthesis of  
 musical notes with wavenet autoencoders. In *Proceedings*  
 of the *International Conference on Machine Learning*,  
 volume 70, pp. 1068–1077. PMLR, 2017.
- European Parliament and Council of the European Union.  
 Regulation (EU) 2016/679 of the European Parliament  
 and of the Council of 27 April 2016 on the protection of  
 natural persons with regard to the processing of personal  
 data and on the free movement of such data, and repealing  
 Directive 95/46/EC (General Data Protection Regulation).  
 EUR-Lex, 2016.
- Fernandes, D., Villa, S., Nicholls, S., Haavisto, O., Buschek,  
 D., Schmidt, A., Kosch, T., Shen, C., and Welsch, R. Ai  
 makes you smarter but none the wiser: The disconnect  
 between performance and metacognition. *Computers in*  
*Human Behavior*, 175:108779, 2026.
- Filippucci, F., Gal, P., Jona-Lasinio, C., Leandro, A., and  
 Nicoletti, G. The impact of Artificial Intelligence on  
 productivity, distribution and growth: Key mechanisms,  
 initial evidence and policy challenges. *OECD Artificial*  
*Intelligence Papers*, (15), 2024.
- Fink, G. An introduction to cultural mismatch theory  
 and its role in equitable learning, April 2023. URL  
[https://www.everylearnereverywhere.org/blog/an-introduction-to-cultural-](https://www.everylearnereverywhere.org/blog/an-introduction-to-cultural-mismatch-theory-and-its-role-in-equitable-learning/)  
[mismatch-theory-and-its-role-in-](https://www.everylearnereverywhere.org/blog/an-introduction-to-cultural-mismatch-theory-and-its-role-in-equitable-learning/)  
[equitable-learning/](https://www.everylearnereverywhere.org/blog/an-introduction-to-cultural-mismatch-theory-and-its-role-in-equitable-learning/). Accessed: 2025-11-06.
- Floridi, L. and COWLS, J. A unified framework of five prin-  
 ciples for AI in society. *Machine learning and the city:*  
*Applications in architecture and urban design*, pp. 535–  
 545, 2022.
- Foerster, H., Behrouzi, S., Rieger, P., Jadhliwala, M., and  
 Sadeghi, A.-R. LightShed: Defeating perturbation-based  
 image copyright protections. In *Proceedings of the 34th*  
*USENIX Security Symposium*, pp. 7271–7290, 2025.
- Frank, J., Herbert, F., Ricker, J., Schönherr, L., Eisenhofer,  
 T., Fischer, A., Dürmuth, M., and Holz, T. A represen-  
 tative study on human detection of artificially generated  
 media across countries. In *2024 IEEE Symposium on*  
*Security and Privacy*, pp. 55–73. IEEE, 2024.
- Gao, R. and Yin, M. Confounding-robust deferral policy  
 learning. In *Proceedings of the AAAI Conference on Arti-*  
*ficial Intelligence*, volume 39, pp. 14238–14246, 2025.
- Gao, W., Liu, Q., Huang, Z., Yin, Y., Bi, H., Wang, M.-C.,  
 Ma, J., Wang, S., and Su, Y. Rcd: Relation map driven  
 cognitive diagnosis for intelligent education systems. In  
*Proceedings of the 44th International ACM SIGIR Con-*  
*ference on Research and Development in Information*  
*Retrieval, SIGIR ’21*, pp. 501–510, New York, NY, USA,  
 2021. Association for Computing Machinery.
- Gatys, L. A., Ecker, A. S., and Bethge, M. A neural algo-  
 rithm of artistic style. arXiv:1508.06576, 2015.

- 605 Gay, G. *Culturally responsive teaching: Theory, research,*  
606 *and practice.* Teachers College Press, 2018.
- 607
- 608 Gerlich, M. AI tools in society: Impacts on cognitive of-  
609 floading and the future of critical thinking. *Societies*, 15  
610 (1), 2025.
- 611
- 612 Ghori, M. F., Dehpanah, A., Gemmell, J., and Mobasher,  
613 B. "they only offer the illusion of choice": Exploring  
614 user perceptions of control and agency on youtube. In  
615 *Adjunct Proceedings of the 33rd ACM Conference on*  
616 *User Modeling, Adaptation and Personalization*, UMAP  
617 Adjunct '25, pp. 214–218, New York, NY, USA, 2025.  
618 Association for Computing Machinery.
- 619
- 620 Gillani, N., Yuan, A., Saveski, M., Vosoughi, S., and Roy,  
621 D. Me, my echo chamber, and I: Introspection on social  
622 media polarization. In *Proceedings of the 2018 World*  
623 *Wide Web Conference*, pp. 823–831, 2018.
- 624
- 625 Gillham, J. Can humans detect AI-generated text on their  
626 own? Forbes Technology Council, 2024.
- 627
- 628 Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B.,  
629 Warde-Farley, D., Ozair, S., Courville, A., and Bengio,  
630 Y. Generative adversarial nets. In *Advances in Neural*  
631 *Information Processing Systems*, volume 27, 2014.
- 632
- 633 Goodier, M. Revealed: Thousands of uk university  
634 students caught cheating using AI. The Guardian,  
635 2025. URL [https://www.theguardian.com/education/2025/jun/15/thousands-](https://www.theguardian.com/education/2025/jun/15/thousands-of-uk-university-students-caught-cheating-using-ai-artificial-intelligence-survey)  
636 [of-uk-university-students-caught-](https://www.theguardian.com/education/2025/jun/15/thousands-of-uk-university-students-caught-cheating-using-ai-artificial-intelligence-survey)  
637 [cheating-using-ai-artificial-](https://www.theguardian.com/education/2025/jun/15/thousands-of-uk-university-students-caught-cheating-using-ai-artificial-intelligence-survey)  
638 [intelligence-survey](https://www.theguardian.com/education/2025/jun/15/thousands-of-uk-university-students-caught-cheating-using-ai-artificial-intelligence-survey). Accessed: 2025-11-  
639 15.
- 640
- 641 Government Offices of Sweden. Chronological overview  
642 of LGBT persons rights in Sweden, 2018. URL  
643 [https://www.government.se/articles/](https://www.government.se/articles/2018/07/chronological-overview-of-lgbt-persons-rights-in-sweden/)  
644 [2018/07/chronological-overview-of-](https://www.government.se/articles/2018/07/chronological-overview-of-lgbt-persons-rights-in-sweden/)  
645 [lgbt-persons-rights-in-sweden/](https://www.government.se/articles/2018/07/chronological-overview-of-lgbt-persons-rights-in-sweden/). Ac-  
646 cessed: 2026-01-03.
- 647
- 648 Griffith, V. and Koch, C. Quantifying synergistic mutual  
649 information. arXiv:1205.4265, 2014.
- 650
- 651 Grinschgl, S., Papenmeier, F., and Meyerhoff, H. S. Conse-  
652 quences of cognitive offloading: Boosting performance  
653 but diminishing memory. *Quarterly Journal of Experi-*  
654 *mental Psychology*, 74(9):1477–1496, 2021.
- 655
- 656 Groh, M., Sankaranarayanan, A., Singh, N., Kim, D. Y.,  
657 Lippman, A., and Picard, R. Human detection of political  
658 speech deepfakes across transcripts, audio, and video.  
659 *Nature Communications*, 15(1):7629, 2024.
- Groot, T. and Valdenegro Toro, M. Overconfidence is key:  
Verbalized uncertainty evaluation in large language and  
vision-language models. In *Proceedings of the 4th Work-*  
*shop on Trustworthy Natural Language Processing*, pp.  
145–171, 2024.
- Gross, R. *Psychology: The Science of Mind and Behaviour*  
*6th Edition.* Hodder Education, 2012.
- Hacker, P. A legal framework for AI training data—from  
first principles to the Artificial Intelligence Act. *Law,*  
*innovation and technology*, 13(2):257–301, 2021.
- Hagendorff, T. The ethics of AI ethics: An evaluation of  
guidelines. *Minds and Machines*, 30(1):99–120, Mar  
2020.
- Hassan, M., Kushniruk, A., and Borycki, E. Barriers to  
and facilitators of artificial intelligence adoption in health  
care: Scoping review. *JMIR Hum Factors*, 11:e48633,  
Aug 2024.
- Heymans, M. Guided learning in Gemini: From an-  
swers to understanding, August 2025. URL [https://blog.google/products-and-platforms/](https://blog.google/products-and-platforms/products/education/guided-learning/)  
[products/education/guided-learning/](https://blog.google/products-and-platforms/products/education/guided-learning/).  
Accessed: 2026-01-28.
- Hirvonen, N., Jylhä, V., Lao, Y., and Larsson, S. Artificial  
intelligence in the information ecosystem: Affordances  
for everyday information seeking. *Journal of the Asso-*  
*ciation for Information Science and Technology*, 75(10):  
1152–1165, 2024.
- Ince, R. A. A. Measuring multivariate redundant informa-  
tion with pointwise common change in surprisal. *Entropy*,  
19(7):318, 2017.
- Jamali, L. AI firm anthropic agrees to pay authors \$1.5bn to  
settle piracy lawsuit, 2025. URL <https://www.bbc.com/news/articles/c5y4jjpg922qo>. Accessed:  
2025-11-15.
- Jiang, B., Hao, Z., Cho, Y.-M., Li, B., Yuan, Y., Chen, S.,  
Ungar, L., Taylor, C. J., and Roth, D. Know me, respond  
to me: Benchmarking LLMs for dynamic user profiling  
and personalized responses at scale. arXiv:2504.14225,  
2025.
- Jobin, A., Ienca, M., and Vayena, E. The global landscape  
of AI ethics guidelines. *Nature machine intelligence*, 1  
(9):389–399, 2019.
- Kambur, H. and Dolunay, A. A research on copyright issues  
impacting artists emotional states in the framework of ar-  
tificial intelligence. *Frontiers in Psychology*, 15:1409646,  
2024.

- 660 Kanervisto, A., Bignell, D., Wen, L. Y., Grayson, M.,  
661 Georgescu, R., Valcarcel Macua, S., Tan, S. Z., Rashid, T.,  
662 Pearce, T., Cao, Y., Lemkhenter, A., Jiang, C., Costello,  
663 G., Gupta, G., Tot, M., Ishida, S., Gupta, T., Arora, U.,  
664 White, R. W., Devlin, S., Morrison, C., and Hofmann,  
665 K. World and human action models towards gameplay  
666 ideation. *Nature*, 638(8051):656–663, February 2025.
- 667 Karras, T., Laine, S., and Aila, T. A style-based generator  
668 architecture for generative adversarial networks. In *Pro-*  
669 *ceedings of the IEEE/CVF conference on computer vision*  
670 *and pattern recognition*, pp. 4401–4410, 2019.
- 671 Keswani, V., Lease, M., and Kenthapadi, K. Towards un-  
672 biased and accurate deferral to multiple experts. In *Pro-*  
673 *ceedings of the AAAI/ACM Conference on AI, Ethics, and*  
674 *Society*, pp. 154–165, 2021.
- 675 Khan, I. Midjourney. In *The Quick Guide to Prompt En-*  
676 *gineering: Generative AI Tips and Tricks for ChatGPT,*  
677 *Bard, Dall-E, and Midjourney*, pp. 333–343. Wiley AI,  
678 2024a.
- 679 Khan, S. *Brave New Words: How AI Will Revolutionize*  
680 *Education (and Why That’s a Good Thing)*. Penguin  
681 Publishing Group, 2024b.
- 682 Kilitcioglu, D., Greenquist, N., and Bari, A. Pyrorank: A  
683 novel nature-inspired algorithm to promote diversity in  
684 recommender systems. In *International Conference on*  
685 *Swarm Intelligence*, pp. 139–155. Springer, 2023.
- 686 Kim, M., Kim, S., Lee, S., Yoon, Y., Myung, J., Yoo, H.,  
687 Lim, H., Han, J., Kim, Y., Ahn, S.-Y., Kim, J., Oh, A.,  
688 Hong, H., and Lee, T. Y. Designing prompt analytics  
689 dashboards to analyze student-chatgpt interactions in efl  
690 writing. arXiv:2405.19691, 2024. URL [https://](https://arxiv.org/abs/2405.19691)  
691 [arxiv.org/abs/2405.19691](https://arxiv.org/abs/2405.19691).
- 692 Kizilcec, R. F. How much information? Effects of trans-  
693 parency on trust in an algorithmic interface. In *Proce-*  
694 *edings of the 2016 CHI Conference on Human Factors in*  
695 *Computing Systems*, pp. 2390–2395, 2016.
- 696 Klock, A. H. and Jan, C. B. Syntax error correction method  
697 and apparatus. [https://patents.google.com/](https://patents.google.com/patent/US4617643A/en)  
698 [patent/US4617643A/en](https://patents.google.com/patent/US4617643A/en), 1986.
- 699 Koh, P. W. and Liang, P. Understanding black-box pre-  
700 dictions via influence functions. In *Proceedings of the*  
701 *International Conference on Machine Learning*, pp. 1885–  
702 1894. PMLR, 2017.
- 703 Kosmyna, N., Hauptmann, E., Yuan, Y. T., Situ, J., Liao,  
704 X.-H., Beresnitzky, A. V., Braunstein, I., and Maes,  
705 P. Your brain on ChatGPT: Accumulation of cognitive  
706 debt when using an AI assistant for essay writing task.  
707 arXiv:2506.08872, 2025.
- 708 Kozyreva, A., Lorenz-Spreen, P., Hertwig, R.,  
709 Lewandowsky, S., and Herzog, S. M. Public atti-  
710 tudes towards algorithmic personalization and use of  
711 personal data online: Evidence from Germany, Great  
712 Britain, and the United States. *Humanities and Social*  
713 *Sciences Communications*, 8(1):1–11, 2021.
- 714 Kuang, W., Qian, B., Li, Z., Chen, D., Gao, D., Pan, X., Xie,  
Y., Li, Y., Ding, B., and Zhou, J. FederatedScope-LLM:  
A comprehensive package for fine-tuning large language  
models in federated learning. In *Proceedings of the ACM*  
*SIGKDD Conference on Knowledge Discovery and Data*  
*Mining*, pp. 5260–5271, 2024.
- Lambert, N., Morrison, J., Pyatkin, V., Huang, S., Ivison,  
H., Brahman, F., Miranda, L. J. V., Liu, A., Dziri, N.,  
Lyu, S., et al. Tulu 3: Pushing frontiers in open language  
model post-training. arXiv:2411.15124, 2024.
- Lau, G. K. R., Niu, X., Dao, H., Chen, J., Foo, C.-S., and  
Low, B. K. H. Waterfall: Scalable framework for robust  
text watermarking and provenance for LLMs. In *Proce-*  
*edings of the 2024 Conference on Empirical Methods in*  
*Natural Language Processing*, pp. 20432–20466, 2024.
- Lau, G. K. R., Dao, H., Lin, N. K. H., and Low, B. K. H.  
Uncertainty quantification for MLLMs. In *Proceedings*  
*of the ICML 2025 Workshop on Reliable and Responsible*  
*Foundation Models*, 2025.
- LeCun, Y. A path towards autonomous machine intelligence.  
*Open Review*, 62(1):1–62, 2022.
- Lee, S., Lai, B., Ryan, F., Boote, B., and Rehg, J. M. Mod-  
eling multimodal social interactions: New challenges  
and baselines with densely aligned representations. In  
*Proceedings of the IEEE/CVF Conference on Computer*  
*Vision and Pattern Recognition*, pp. 14585–14595, 2024.
- Leopold, T. How AI is reshaping the career lad-  
der, and other trends in jobs and skills on labour  
day. World Economic Forum, 2025. URL [https://](https://www.weforum.org/stories/2025/04/ai-jobs-international-workers-day/)  
[www.weforum.org/stories/2025/04/ai-](https://www.weforum.org/stories/2025/04/ai-jobs-international-workers-day/)  
[jobs-international-workers-day/](https://www.weforum.org/stories/2025/04/ai-jobs-international-workers-day/). Ac-  
cessed: 2026-01-28.
- Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V.,  
Goyal, N., Küttler, H., Lewis, M., Yih, W.-t., Rocktäschel,  
T., Riedel, S., and Kiela, D. Retrieval-augmented gener-  
ation for knowledge-intensive NLP tasks. In *Advances in*  
*Neural Information Processing Systems*, volume 33, pp.  
9459–9474, 2020.
- Li, C., Chen, M., Wang, J., Sitaram, S., and Xie, X. Cul-  
tureLLM: Incorporating cultural differences into large  
language models. In *Advances in Neural Information*  
*Processing Systems*, volume 37, pp. 84799–84838, 2024a.

- 715 Li, F. *The Worlds I See: Curiosity, Exploration, and Discovery at the Dawn of AI*. Flatiron Books, 2025. ISBN  
716 9781250898104. URL <https://books.google.com.sg/books?id=T6bvEAAAQBAJ>.  
717  
718
- 719 Li, L. H., Zhang, P., Zhang, H., Yang, J., Li, C., Zhong, Y.,  
720 Wang, L., Yuan, L., Zhang, L., Hwang, J.-N., Chang, K.-  
721 W., and Gao, J. Grounded language-image pre-training.  
722 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10965–10975,  
723 2022.  
724  
725
- 726 Li, R., Zhang, S., Lin, D., Chen, K., and He, X. From pixels  
727 to graphs: Open-vocabulary scene graph generation with  
728 vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,  
729 pp. 28076–28086, 2024b.  
730  
731
- 732 Li, Y., Xiong, M., Wu, J., and Hooi, B. Conftuner: Training  
733 large language models to express their confidence verbally.  
734 In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025.  
735
- 736 Liang, P. P., Wu, C., Morency, L.-P., and Salakhutdinov,  
737 R. Towards understanding and mitigating social biases in  
738 language models. In *International conference on machine learning*, pp. 6565–6576. PMLR, 2021.  
739  
740
- 741 Liang, P. P., Cheng, Y., Fan, X., Ling, C. K., Nie, S., Chen,  
742 R. J., Deng, Z., Allen, N., Auerbach, R., Mahmood, F.,  
743 Salakhutdinov, R., and Morency, L.-P. Quantifying &  
744 modeling multimodal interactions: An information decomposition framework. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.  
745  
746
- 747 Liu, B., Schlegel, V., Batista-Navarro, R., and Ananiadou,  
748 S. Argument mining as a multi-hop generative machine  
749 reading comprehension task. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp.  
750 10846–10858, 2023.  
751  
752
- 753 Liu, J., Huang, Z., Xiao, T., Sha, J., Wu, J., Liu, Q., Wang,  
754 S., and Chen, E. Socraticlm: Exploring socratic personalized teaching with large language models. In Globerson,  
755 A., Mackey, L., Belgrave, D., Fan, A., Paquet, U., Tomczak, J., and Zhang, C. (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 85693–85721.  
756  
757 Curran Associates, Inc., 2024a.  
758  
759
- 760 Liu, Q., Zhou, Y., Huang, J., and Li, G. When ChatGPT  
761 is gone: Creativity reverts and homogeneity persists.  
762 arXiv:2401.06816, 2024b.  
763  
764
- 765 Liu, S., Shen, J., Qian, H., and Zhou, A. Inductive cognitive diagnosis for fast student learning in web-based  
766 intelligent education systems. In *Proceedings of the ACM Web Conference 2024, WWW '24*, pp. 4260–4271. ACM,  
767 May 2024c.  
768  
769
- Liu, S., Zeng, Z., Ren, T., Li, F., Zhang, H., Yang, J., Jiang, Q., Li, C., Yang, J., Su, H., Zhu, J., and Zhang, L. Grounding DINO: Marrying DINO with grounded pre-training for open-set object detection. In *Proceedings of the European Conference on Computer Vision*, pp. 38–55, 2024d.
- Liu, Y., Fan, C., Dai, Y., Chen, X., Zhou, P., and Sun, L. Metacloak: Preventing unauthorized subject-driven text-to-image diffusion-based synthesis via meta-learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 24219–24228, 2024e.
- Liu, Z., Lin, G., Tan, H. L., Zhang, H., Lu, Y., Gao, X., Yin, S. X., He, S., Goh, H. H., Wong, L. H., and Chen, N. F. SingaKids: A multilingual multimodal dialogic tutor for language learning. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 6: Industry Track)*, pp. 1244–1253, 2025.
- Lovato, J., Zimmerman, J., Smith, I., Dodds, P., and Karson, J. Foregrounding artist opinions: A survey study on transparency, ownership, and fairness in AI generative art. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, volume 7, pp. 905–916, 2024.
- Lu, X., Niu, X., Lau, G. K. R., Nhung, B. T. C., Sim, R. H. L., Wen, F., Foo, C.-S., Ng, S.-K., and Low, B. K. H. Waterdrum: Watermarking for data-centric unlearning metric. arXiv:2505.05064, 2025a.
- Lu, X., Sclar, M., Hallinan, S., Miresghallah, N., Liu, J., Han, S., Ettinger, A., Jiang, L., Chandu, K., Dziri, N., and Choi, Y. AI as humanity’s salieri: Quantifying linguistic creativity of language models via systematic attribution of machine text against web text. In *Proceedings of the International Conference on Learning Representations*, 2025b.
- Lu, X., Wang, J., Zhao, Z., Dai, Z., Foo, C.-S., Ng, S.-K., and Low, B. K. H. WASA: Watermark-based Source Attribution for large language model-generated data. In *Findings of the Association for Computational Linguistics*, pp. 23791–23824, 2025c.
- Luccioni, S., Akiki, C., Mitchell, M., and Jernite, Y. Stable bias: Evaluating societal representations in diffusion models. In *Advances in Neural Information Processing Systems*, volume 36, pp. 56338–56351, 2023.
- Ma, P., Yang, X., Li, Y., Gui, M., Krause, F., Schusterbauer, J., and Ommer, B. SCFlow: Implicitly learning style and content disentanglement with flow models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 14919–14929, 2025.
- Macnamara, B. N., Berber, I., Çavuşoğlu, M. C., Krupinski, E. A., Nallapareddy, N., Nelson, N. E., Smith, P. J.,

- 770 Wilson-Delfosse, A. L., and Ray, S. Does using artificial  
771 intelligence assistance accelerate skill decay and hinder  
772 skill development without performers’ awareness? *Cog-  
773 nitive Research: Principles and Implications*, 9(1):46,  
774 2024.
- 775  
776 Magesh, V., Surani, F., Dahl, M., Suzgun, M., Manning,  
777 C. D., and Ho, D. E. Hallucination-free? Assessing the  
778 reliability of leading AI legal research tools. *Journal of  
779 Empirical Legal Studies*, 22(2):216–242, 2025.
- 780  
781 Maitland, E. and Lee, A. 400 and counting: A Russian influ-  
782 ence operation overtakes official state media in spreading  
783 Russia-Ukraine false claims, 2025. URL <https://www.newsguardtech.com/special-reports/400-and-counting-a-russian-influence-operation-overtakes-official-state-media-in-spreading-russia-ukraine-false-claims/>. Accessed:  
784 2026-01-28.
- 785  
786  
787  
788  
789  
790 Martela, F. The normative value of making a positive  
791 contribution–benefiting others as a core dimension of  
792 meaningful work. *Journal of Business Ethics*, 185(4):  
793 811–823, 2023.
- 794  
795 Masters, C., Vellanki, A., Shangguan, J., Kultys, B.,  
796 Gilmore, J., Moore, A., and Albrecht, S. Orchestrating  
797 human-AI teams: The manager agent as a unifying  
798 research challenge. In *Proceedings of the International  
799 Conference on Distributed Artificial Intelligence*, pp. 91–  
800 107, 2025.
- 801  
802 Matthias, A. The responsibility gap: Ascribing responsi-  
803 bility for the actions of learning automata. *Ethics and  
804 information technology*, 6(3):175–183, 2004.
- 805  
806 Maurya, K. K., Srivatsa, K. A., Petukhova, K., and Kochmar,  
807 E. Unifying AI tutor evaluation: An evaluation taxonomy  
808 for pedagogical ability assessment of LLM-powered ai  
809 tutors. In *Proceedings of the 2025 Conference of the  
810 Nations of the Americas Chapter of the Association for  
811 Computational Linguistics: Human Language Technolo-  
812 gies (Volume 1: Long Papers)*, pp. 1234–1251, 2025.
- 813  
814 McMahan, B. J., Peng, Z., Zhou, B., and Kao, J. C. Shared  
815 autonomy with IDA: interventional diffusion assistance.  
816 In *Advances in Neural Information Processing Systems*,  
817 volume 37, pp. 128330–128354, 2024.
- 818  
819 Mentutor. Effortless education: Hosting, learning,  
820 and teaching with mentutor, 2024. URL <https://www.mentutor.io/post/effortless-education-hosting-learning-and-teaching-with-mentutor>. Accessed: 2025-12-  
821 20.
- 822  
823  
824
- Merchant, B. AI is already taking jobs in the video  
game industry, 2023. URL <https://www.wired.com/story/ai-is-already-taking-jobs-in-the-video-game-industry/>. Accessed:  
2025-11-09.
- Meyers, S., Rowell, K., Wells, M., and Smith, B. C. Teacher  
empathy: A model of empathy for teaching for student  
success. *College Teaching*, 67(3):160–168, July 2019.
- Mills, S. and Sætra, H. S. The autonomous choice architect.  
*AI & Society*, 39(2):583–595, 2024.
- Min, S., Gururangan, S., Wallace, E., Shi, W., Hajishirzi, H.,  
Smith, N. A., and Zettlemoyer, L. SILO language mod-  
els: Isolating legal risk in a nonparametric datastore. In  
*Proceedings of the International Conference on Learning  
Representations*, 2024.
- Mlodozieniec, B. K., Eschenhagen, R., Bae, J., Immer, A.,  
Krueger, D., and Turner, R. E. Influence functions for  
scalable data attribution in diffusion models. In *Pro-  
ceedings of the International Conference on Learning  
Representations*, 2025.
- Montgomery, B. Disney and Universal sue AI im-  
age creator Midjourney, alleging copyright infringe-  
ment, 2025. URL <https://www.theguardian.com/technology/2025/jun/11/disney-universal-ai-lawsuit>. Accessed: 2025-12-29.
- Moon, K., Green, A. E., and Kushlev, K. Homogenizing  
effect of large language models (LLMs) on creative di-  
versity: An empirical comparison of human and chatgpt  
writing. *Computers in Human Behavior: Artificial Hu-  
mans*, 6:100207, 2025.
- Morley, J., Elhalal, A., Garcia, F., Kinsey, L., Mökander, J.,  
and Floridi, L. Ethics as a service: A pragmatic opera-  
tionalisation of AI ethics. *Minds and Machines*, 31(2):  
239–256, 2021.
- Morley, J., Kinsey, L., Elhalal, A., Garcia, F., Ziosi, M., and  
Floridi, L. Operationalising AI ethics: barriers, enablers  
and next steps. *AI & Society*, 38(1):411–423, 2023.
- Naous, T., Ryan, M. J., Ritter, A., and Xu, W. Having beer  
after prayer? Measuring cultural bias in large language  
models. In *Proceedings of the 62nd Annual Meeting of  
the Association for Computational Linguistics (Volume 1:  
Long Papers)*, pp. 16366–16393, 2024.
- Natali, C., Marconi, L., Dias Duran, L. D., and Cabitza,  
F. AI-induced deskilling in medicine: A mixed-method  
review and research agenda for healthcare and beyond.  
*Artificial Intelligence Review*, 58(11), 2025.

- 825 Ng, R., Nguyen, T. N., Huang, Y., Tai, N. C., Leong, W. Y.,  
826 Leong, W. Q., Yong, X., Ngui, J. G., Susanto, Y., Cheng,  
827 N., Rengarajan, H., Limkonchotiwat, P., Hulagadri, A. V.,  
828 Teng, K. W., Tong, Y. Y., Siow, B., Teo, W. Y., Lau, W.,  
829 Tan, C. M., Ong, B., Ong, Z. H., Montalan, J. R., Chan,  
830 A., Antonyrex, S., Lee, R., Choa, E., Tat-Wee, D. O., Liu,  
831 B. J. D., Tjhi, W. C., Cambria, E., and Teo, L. SEA-LION:  
832 Southeast asian languages in one network. In *Proceedings*  
833 *of the 14th International Joint Conference on Natural*  
834 *Language Processing and the 4th Conference of the Asia-*  
835 *Pacific Chapter of the Association for Computational*  
836 *Linguistics*, pp. 512–526, 2025.
- 837  
838 Ng, S. F. The model method: Crown jewel in singapore  
839 mathematics. *Asian Journal for Mathematics Education*,  
840 1(2):147–161, 2022.
- 841  
842 Ngo, T., Kunkel, J., and Ziegler, J. Exploring mental models  
843 for transparent and controllable recommender systems: A  
844 qualitative study. In *Proceedings of the 28th ACM Confer-*  
845 *ence on User Modeling, Adaptation and Personalization*,  
846 pp. 183–191, 2020.
- 847  
848 Nicoletti, L. and Bass, D. Humans are biased. Gener-  
849 ative AI is even worse, 2023. URL <https://www.bloomberg.com/graphics/2023-generative-ai-bias/>. Accessed: 2026-01-28.
- 850  
851  
852 Nikolic, I., Baluta, T., and Saxena, P. Model provenance  
853 testing for large language models. arXiv:2502.00706,  
854 2025.
- 855  
856 Ning, L., Liu, L., Wu, J., Wu, N., Berlowitz, D., Prakash,  
857 S., Green, B., O’Banion, S., and Xie, J. User-LLM: Effi-  
858 cient LLM contextualization with user embeddings. In  
859 *Companion Proceedings of the ACM on Web Conference*  
860 *2025*, pp. 1219–1223, 2025.
- 861  
862 Nolan, B. AI is gutting the next generation of talent: In tech,  
863 job openings for new grads have already been halved,  
864 2025. URL <https://fortune.com/2025/08/15/ai-gutting-next-generation-of-talent/>. Accessed: 2025-12-29.
- 865  
866  
867  
868 Noy, S. and Zhang, W. Experimental evidence on the produc-  
869 tivity effects of generative artificial intelligence. *Science*,  
870 381(6654):187–192, 2023.
- 871  
872  
873 Oh, S.-y. and Ahn, Y. Exploring teachers’ perception of  
874 artificial intelligence: The socio-emotional deficiency as  
875 opportunities and challenges in human-ai complementar-  
876 ity in K-12 education. In Olney, A. M., Chounta, I.-A.,  
877 Liu, Z., Santos, O. C., and Bittencourt, I. I. (eds.), *Artifi-*  
878 *cial Intelligence in Education*, pp. 439–447, Cham, 2024.  
879 Springer Nature Switzerland.
- OpenAI. Introducing study mode, 2025. URL <https://openai.com/index/chatgpt-study-mode/>. Accessed: 2026-01-28.
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C. L., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., Schulman, J., Hilton, J., Kelton, F., Miller, L., Simens, M., Askell, A., Welinder, P., Christiano, P., Leike, J., and Lowe, R. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems*, volume 35, pp. 27730–27744, 2022.
- Paivio, A., Rogers, T. B., and Smythe, P. C. Why are pictures easier to recall than words? *Psychonomic Science*, 11(4): 137–138, Apr 1968.
- Papakyriakopoulos, O. and Goodman, E. The impact of twitter labels on misinformation spread and user engagement: Lessons from trump’s election tweets. In *Proceedings of the ACM Web Conference 2022*, pp. 2541–2551, 2022.
- Park, S. The work of art in the age of generative AI: aura, liberation, and democratization. *AI & society*, 40(3): 1807–1816, 2025.
- Paul, C. and Matthews, M. The Russian “firehose of falsehood” propaganda model: Why it might work and options to counter it. *Rand Corporation*, 2(7):1–10, 2016.
- Peper, J., Qiu, W., and Wang, L. PELMS: Pre-training for effective low-shot multi-document summarization. In Duh, K., Gomez, H., and Bethard, S. (eds.), *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 7652–7674, Mexico City, Mexico, June 2024. Association for Computational Linguistics.
- Pope, A. NYT v. OpenAI: The Times’s about-face, 2024. URL <https://harvardlawreview.org/blog/2024/04/nyt-v-openai-the-timess-about-face/>. Accessed: 2026-01-17.
- Puig, X., Shu, T., Li, S., Wang, Z., Liao, Y.-H., Tenenbaum, J. B., Fidler, S., and Torralba, A. Watch-And-Help: A challenge for social perception and human-AI collaboration. In *Proceedings of the International Conference on Learning Representations*, 2021.
- Qi, N., Li, Y., Fu, R., and Zhu, Q. STRAT: Image style transfer with region-aware transformer. *Neurocomputing*, 617:129039, 2025.
- Qi, Y., Schölkopf, B., and Jin, Z. Causal responsibility attribution for human-AI collaboration. arXiv:2411.03275, 2024.

- 880 Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G.,  
 881 Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark,  
 882 J., Krueger, G., and Sutskever, I. Learning transferable  
 883 visual models from natural language supervision. In  
 884 *Proceedings of the International Conference on Machine*  
 885 *Learning*, volume 139, pp. 8748–8763. PMLR, 2021.
- 886 Rafailov, R., Sharma, A., Mitchell, E., Ermon, S., Manning,  
 887 C. D., and Finn, C. Direct preference optimization: your  
 888 language model is secretly a reward model. In *Advances*  
 889 *in Neural Information Processing Systems*, volume 36,  
 890 pp. 53728–53741, 2023.
- 891 Raji, I. D., Smart, A., White, R. N., Mitchell, M., Gebru, T.,  
 892 Hutchinson, B., Smith-Loud, J., Theron, D., and Barnes,  
 893 P. Closing the AI accountability gap: defining an end-  
 894 to-end framework for internal algorithmic auditing. In  
 895 *Proceedings of the 2020 Conference on Fairness, Ac-*  
 896 *countability, and Transparency*, pp. 33–44. Association  
 897 for Computing Machinery, 2020.
- 898 Reuters. Chinese court rejects landmark same-  
 899 sex marriage case, 2016. URL <https://www.reuters.com/article/world/chinese-court-rejects-landmark-same-sex-marriage-case-idUSKCN0XA1BR/>.  
 900 Accessed: 2026-01-03.
- 901 Risko, E. F. and Gilbert, S. J. Cognitive offloading. *Trends*  
 902 *in cognitive sciences*, 20(9):676–688, 2016.
- 903 Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and  
 904 Ommer, B. High-resolution image synthesis with la-  
 905 tent diffusion models. In *Proceedings of the IEEE/CVF*  
 906 *Conference on Computer Vision and Pattern Recognition*  
 907 (CVPR), pp. 10684–10695, 2022.
- 908 Rouinfar, A., Agra, E., Larson, A., Loschky, L., and Rebello,  
 909 N. Can visual cues and correctness feedback influence  
 910 students’ reasoning. In *American Institute of Physics*  
 911 *Conference Series*, 2014.
- 912 Rowan, L. and Grootenboer, P. (eds.). *Student engagement*  
 913 *and educational rapport in higher education*. Springer  
 914 International Publishing, Cham, Switzerland, 1 edition,  
 915 November 2016.
- 916 Ruiz-Dolz, R., Kikteva, Z., and Lawrence, J. Mining com-  
 917 plex patterns of argumentative reasoning in natural lan-  
 918 guage dialogue. In *Proceedings of the 63rd Annual Meet-*  
 919 *ing of the Association for Computational Linguistics (Vol-*  
 920 *ume 1: Long Papers)*, pp. 7421–7435, 2025.
- 921 Ryu, S., Do, H., Kim, Y., Lee, G., and Ok, J. Multi-  
 922 dimensional optimization for text summarization via rein-  
 923 forcement learning. In Ku, L.-W., Martins, A., and Sriku-  
 924 mar, V. (eds.), *Proceedings of the 62nd Annual Meeting*  
 925 *of the Association for Computational Linguistics (Volume*  
 926 *1: Long Papers)*, pp. 5858–5871, Bangkok, Thailand,  
 927 August 2024. Association for Computational Linguistics.
- 928 Sackett, P. R. and Walmsley, P. T. Which personality at-  
 929 tributes are most important in the workplace? *Perspec-*  
 930 *tives on Psychological Science*, 9(5):538–551, 2014.
- 931 Sanchini, V., Pongiglione, F., and Sala, R. On the notion of  
 932 political agency. *Phenomenology and Mind*, 16:10–15,  
 933 2019.
- 934 Seah, M. Statement by Mr Mark Seah, Deputy Perma-  
 935 nent Representative of Singapore, on Agenda Item 84,  
 936 on the rule of law at the national and international levels,  
 937 sixth committee. Permanent Mission of Singapore to the  
 938 United Nations, 2022. URL [https://www.un.org/en/ga/sixth/77/pdfs/statements/rule\\_of\\_law/06mtg\\_singapore.pdf](https://www.un.org/en/ga/sixth/77/pdfs/statements/rule_of_law/06mtg_singapore.pdf). Accessed: 2025-12-29.
- 939 Senftleben, M. Generative AI and author remuneration.  
 940 *IIC-International Review of Intellectual Property and*  
 941 *Competition Law*, 54(10):1535–1560, 2023.
- 942 Seo, S., Han, B., Harari, R. E., Dias, R. D., Zenati, M. A.,  
 943 Salas, E., and Unhelkar, V. Socratic: Enhancing human  
 944 teamwork via AI-enabled coaching. In *Proceedings of*  
 945 *the International Conference on Autonomous Agents and*  
 946 *Multiagent Systems*, pp. 1876–1885, 2025.
- 947 Shabanov, I. Unlock your academic potential. Get the ef-  
 948 fortless academic quick start guide, 2025. URL <https://effortlessacademic.com/>. Accessed: 2025-12-29.
- 949 Shah, M. and Sureja, N. A comprehensive review of bias  
 950 in deep learning models: Methods, impacts, and future  
 951 directions. *Archives of Computational Methods in Engi-*  
 952 *neering*, 32(1):255–267, 2025.
- 953 Shamshad, F., Bakr, T., Shaaban, Y., Hussein, N., Nan-  
 954 dakumar, K., and Lukas, N. First-place solution to  
 955 NeurIPS 2024 invisible watermark removal challenge.  
 956 arXiv:2508.21072, 2025.
- 957 Shan, S., Cryan, J., Wenger, E., Zheng, H., Hanocka, R.,  
 958 and Zhao, B. Y. Glaze: protecting artists from style  
 959 mimicry by text-to-image models. In *Proceedings of the*  
 960 *32nd USENIX Conference on Security Symposium*, pp.  
 961 2187–2204, 2023.
- 962 Shan, S., Ding, W., Passananti, J., Wu, S., Zheng, H., and  
 963 Zhao, B. Y. Nightshade: Prompt-specific poisoning at-  
 964 tacks on text-to-image generative models. In *2024 IEEE*  
 965 *Symposium on Security and Privacy*, pp. 807–825. IEEE,  
 966 2024.

- 935 Shanmugasundaram, M. and Tamilarasu, A. The impact of  
936 digital technology, social media, and artificial intelligence  
937 on cognitive functions: a review. *Frontiers in Cognition*,  
938 2(1203077), 2023.
- 939 Shi, H., Xu, Z., Wang, H., Qin, W., Wang, W., Wang, Y.,  
940 Wang, Z., Ebrahimi, S., and Wang, H. Continual learning  
941 of large language models: A comprehensive survey. *ACM*  
942 *Computing Surveys*, 58(5):1–42, 2025.
- 944 Shu, Y., Hu, W., Ng, S.-K., Low, B. K. H., and Yu, F. Ferret:  
945 Federated full-parameter tuning at scale for large lan-  
946 guage models. In *Proceedings of the International Con-*  
947 *ference on Machine Learning*, pp. 55379–55402. PMLR,  
948 2025.
- 949 Shuai, Z., Wu, C., Tang, Z., Song, B., and Shen, L. Latent  
950 space disentanglement in diffusion transformers enables  
951 precise zero-shot semantic editing. arXiv:2411.08196,  
952 2024.
- 954 Shumailov, I., Shumaylov, Z., Zhao, Y., Papernot, N., Ander-  
955 son, R., and Gal, Y. AI models collapse when trained on  
956 recursively generated data. *Nature*, 631(8022):755–759,  
957 2024.
- 958 Sikander, G. and Anwar, S. Driver fatigue detection systems:  
959 A review. *IEEE Transactions on Intelligent Transporta-*  
960 *tion Systems*, 20(6):2339–2352, 2018.
- 962 Simon, L. K. Is AI responsible for the rise in entry-level  
963 unemployment? Revelop Labs, 2025. URL <https://www.reveliolabs.com/news/macro/is-ai-responsible-for-the-rise-in-entry-level-unemployment/>.
- 964 Skibba, R. Bosses say AI boosts productivity – workers  
965 say they’re drowning in ‘workslop’. The Guardian,  
966 2026. URL <https://www.theguardian.com/technology/2026/apr/14/ai-productivity-workplace-errors>. Accessed:  
967 2026-05-07.
- 968 Slokom, M., Hanjalic, A., and Larson, M. Towards  
969 user-oriented privacy for recommender system data: A  
970 personalization-based approach to gender obfuscation for  
971 user profiles. *Information Processing & Management*, 58  
972 (6):102722, 2021.
- 973 Solovev, K. and Pröllochs, N. References to unbiased  
974 sources increase the helpfulness of community fact-  
975 checks. *Scientific Reports*, 15(1):25749, 2025.
- 976 South China Morning Post. ChatGPT usage hits record as  
977 Studio Ghibli-style AI images go viral, April 2025. URL  
978 <https://www.scmp.com/tech/big-tech/article/3304828/chatgpt-usage-hits-record-studio-ghibli-style-ai-images-go-viral>. Accessed: 2025-10-30.
- 980 Stadler, M., Bannert, M., and Sailer, M. Cognitive ease at a  
981 cost: Llms reduce mental effort but compromise depth in  
982 student scientific inquiry. *Computers in Human Behavior*,  
983 160:108386, 2024.
- 984 Steele, D. M. and Cohn-Vargas, B. *Identity safe classrooms,*  
985 *grades K-5: Places to belong and learn*. Corwin Press,  
986 2013.
- 987 Sung, Y.-T., Yang, J.-M., and Han, Y. L. The effects of  
988 mobile-computer-supported collaborative learning: Meta-  
989 analysis and critical synthesis. *Review of Educational*  
*Research*, 87(4):768–805, 2017.
- Sutskever, I., Martens, J., and Hinton, G. Generating  
text with recurrent neural networks. In *Proceedings of*  
*the International Conference on Machine Learning*, pp.  
1017–1024, 2011.
- Tan, J., Xu, S., Ge, Y., Li, Y., Chen, X., and Zhang, Y. Counterfactual explainable recommendation. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, pp. 1784–1793, 2021.
- Tan, J., Ge, Y., Zhu, Y., Xia, Y., Luo, J., Ji, J., and Zhang, Y. User-controllable recommendation via counterfactual retrospective and prospective explanations. In *Proceedings of the 26th European Conference on Artificial Intelligence*, pp. 2307–2314, 2023.
- Tan, X., Cheng, G., and Ling, M. H. Artificial intelligence in teaching and teacher professional development: A systematic review. *Computers and Education: Artificial Intelligence*, 8:100355, 2025.
- Teng, D., Ye, C., and Martinez, V. AI is stronger than a new graduate, why will you recruit now? A human capital perspective. 2024.
- The New York Times. Exhibit J: One hundred examples of GPT-4 memorizing content from The New York Times, 2023. URL <https://nytimes.com/2023/12/Lawsuit-Document-dkt-1-68-Ex-J.pdf>. Accessed: 2026-01-17.
- Tian, Q. and Zheng, X. Effectiveness of online collaborative problem-solving method on students’ learning performance: A meta-analysis. *Journal of Computer Assisted Learning*, 40(1):326–341, 2024.
- Tian, Z., Liu, A., Esbenshade, L., Sarkar, S., Zhang, Z., He, K., and Sun, M. Implementation considerations for automated AI grading of student work. In *Proceedings of the Artificial Intelligence in Measurement and Education Conference (AIME-Con): Full Papers*, pp. 9–20, 2025.
- Tik Tok. For you. Tik Tok, 2025. URL <https://support.tiktok.com/en/>

- 990 [getting-started/for-you](#). Accessed: 2025-08-  
991 29.
- 992 Tommasel, A. and Menczer, F. Do recommender systems  
993 make social media more susceptible to misinformation  
994 spreaders? In *Proceedings of the 16th ACM Conference*  
995 *on Recommender Systems*, pp. 550–555, 2022.
- 997 Turnitin. Turnitin similarity: Comprehensive plagiarism  
998 detection, 2025. URL <https://www.turnitin.com/products/similarity/>. Accessed: 2025-  
999 11-06.
- 1000  
1001  
1002 United Nations. Universal declaration of human rights, 1948.  
1003 URL [https://www.un.org/en/about-us/](https://www.un.org/en/about-us/universal-declaration-of-human-rights)  
1004 [universal-declaration-of-human-rights](https://www.un.org/en/about-us/universal-declaration-of-human-rights).  
1005 Accessed: 2026-01-03.
- 1006 United Nations Conference on Trade and Development. *Cre-*  
1007 *ative economy outlook 2024*. United Nations, 2024.
- 1009 United Nations Human Rights Office of the High  
1010 Commissioner. Convention against torture and other  
1011 cruel, inhuman or degrading treatment or punishment,  
1012 1984. URL [https://www.ohchr.org/en/](https://www.ohchr.org/en/instruments-mechanisms/instruments/convention-against-torture-and-other-cruel-inhuman-or-degrading)  
1013 [instruments - mechanisms / instruments /](https://www.ohchr.org/en/instruments-mechanisms/instruments/convention-against-torture-and-other-cruel-inhuman-or-degrading)  
1014 [convention-against-torture-and-other-](https://www.ohchr.org/en/instruments-mechanisms/instruments/convention-against-torture-and-other-cruel-inhuman-or-degrading)  
1015 [cruel - inhuman - or - degrading](https://www.ohchr.org/en/instruments-mechanisms/instruments/convention-against-torture-and-other-cruel-inhuman-or-degrading). Accessed:  
1016 2026-01-03.
- 1017 United Nations Human Rights Office of the High  
1018 Commissioner. Committee on the Rights of  
1019 the Child reviews the report of Singapore, 2019.  
1020 URL [https://www.ohchr.org/en/press-](https://www.ohchr.org/en/press-releases/2019/05/committee-rights-child-reviews-report-singapore)  
1021 [releases / 2019 / 05 / committee - rights -](https://www.ohchr.org/en/press-releases/2019/05/committee-rights-child-reviews-report-singapore)  
1022 [child - reviews - report - singapore](https://www.ohchr.org/en/press-releases/2019/05/committee-rights-child-reviews-report-singapore). Ac-  
1023 cessed: 2026-01-03.
- 1025 U.S. Copyright Office. What is copyright? URL [https://](https://www.copyright.gov/what-is-copyright/)  
1026 [www.copyright.gov/what-is-copyright/](https://www.copyright.gov/what-is-copyright/).  
1027 Accessed: 2025-11-20.
- 1028  
1029 U.S. Mission Geneva. U.S. explanation of vote on  
1030 the right to food. U.S. Mission to International  
1031 Organizations in Geneva, 2017. URL [https://](https://geneva.usmission.gov/2017/03/24/u-s-explanation-of-vote-on-the-right-to-food/)  
1032 [geneva.usmission.gov/2017/03/24/u-s-](https://geneva.usmission.gov/2017/03/24/u-s-explanation-of-vote-on-the-right-to-food/)  
1033 [explanation-of-vote-on-the-right-to-](https://geneva.usmission.gov/2017/03/24/u-s-explanation-of-vote-on-the-right-to-food/)  
1034 [food/](https://geneva.usmission.gov/2017/03/24/u-s-explanation-of-vote-on-the-right-to-food/). Accessed: 2025-11-20.
- 1035  
1036 Vaccaro, M., Almaatouq, A., and Malone, T. When com-  
1037 binations of humans and AI are useful: A systematic  
1038 review and meta-analysis. *Nature Human Behaviour*, 8  
1039 (12):2293–2303, 2024.
- 1040 van den Oord, A., Dieleman, S., Zen, H., Simonyan, K.,  
1041 Vinyals, O., Graves, A., Kalchbrenner, N., Senior, A., and  
1042 Kavukcuoglu, K. WaveNet: A generative model for raw  
1043 audio. arXiv:1609.03499, 2016.
- 1044
- Vie, J.-J., Popineau, F., Bruillard, É., and Bourda, Y. A re-  
view of recent advances in adaptive assessment. *Learning*  
*analytics: Fundaments, applications, and trends: A view*  
*of the current state of the art to enhance e-learning*, pp.  
113–142, 2017.
- Wadinambiarachchi, S., Kelly, R. M., Pareek, S., Zhou, Q.,  
and Velloso, E. The effects of generative AI on design  
fixation and divergent thinking. In *Proceedings of the*  
*2024 CHI Conference on Human Factors in Computing*  
*Systems*, pp. 1–18, 2024.
- Wang, J. and Fan, W. The effect of ChatGPT on students’  
learning performance, learning perception, and higher-  
order thinking: insights from a meta-analysis. *Humanities*  
*and Social Sciences Communications*, 12(1):621, 2025.
- Wang, R., Yu, H., Zhang, W., Qi, Z., Sap, M., Bisk, Y.,  
Neubig, G., and Zhu, H. SOTOPIA- $\pi$ : Interactive learn-  
ing of socially intelligent language agents. In *Proceed-*  
*ings of the 62nd Annual Meeting of the Association for*  
*Computational Linguistics (Volume 1: Long Papers)*, pp.  
12912–12940, 2024a.
- Wang, W., Feng, F., Nie, L., and Chua, T.-S. User-  
controllable recommendation against filter bubbles. In  
*Proceedings of the 45th International ACM SIGIR Con-*  
*ference on Research and Development in Information*  
*Retrieval*, pp. 1251–1261, 2022.
- Wang, W., Lin, X., Wang, L., Feng, F., Ma, Y., and Chua,  
T.-S. Causal disentangled recommendation against user  
preference shifts. *ACM Transactions on Information Sys-*  
*tems*, 42(1):1–27, 2023.
- Wang, Y., Li, Z., Zhang, W., Zhang, Z., Xie, B., Liu, X.,  
Zeng, W., and Jin, X. Scene graph disentanglement and  
composition for generalizable complex image generation.  
In *Advances in Neural Information Processing Systems*,  
volume 37, pp. 98478–98504, 2024b.
- Watiktinnakorn, C., Seesai, J., and Kerdvibulvech, C. Blur-  
ring the lines: how AI is redefining artistic ownership and  
copyright. *Discover Artificial Intelligence*, 3(1):37, 2023.
- Williams, P. L. and Beer, R. D. Nonnegative decomposition  
of multivariate information. arXiv:1004.2515, 2010.
- Wilson, R. A. Afterword to “Anthropology and Human  
Rights in a New Key”: The social life of human rights.  
*American Anthropologist*, 108(1):77–83, 2006.
- Winner, L. Do artifacts have politics? *Daedalus*, 109(1):  
121–136, 1980.
- Woods, D. D. Paradigms for intelligent decision support. In  
*Intelligent Decision Support in Process Environments*, pp.  
153–173, 1986.

- World Intellectual Property Organization. Copyright basics: Ideas and expression. World Intellectual Property Organization Website, 2025. URL <https://www.wipo.int/copyright/en/>. Accessed: 2025-10-15.
- World Trade Organization. Module II: Copyright and related rights. Training Module on the TRIPS Agreement, 1995. URL [https://www.wto.org/english/tratop\\_e/trips\\_e/ta\\_docs\\_e/modules2\\_e.pdf](https://www.wto.org/english/tratop_e/trips_e/ta_docs_e/modules2_e.pdf). Accessed: 2025-10-15.
- Wu, X., Li, Y., Zu, T., Hutson, J., Loschky, L. C., and Rebello, N. S. Using multimedia hints to facilitate conceptual problem solving in physics: Investigating the effects of multiple modalities. In *Frontiers in Education*, volume 10, pp. 1568406. Frontiers Media SA, 2025.
- Yakura, H., Lopez-Lopez, E., Brinkmann, L., Serna, I., Gupta, P., Soraperra, I., and Rahwan, I. Empirical evidence of large language model’s influence on human spoken communication. arXiv:2409.01754, 2025.
- Yang, T., Wang, Y., Lv, Y., and Zheng, N. DisDiff: Unsupervised disentanglement of diffusion probabilistic models. arXiv:2301.13721, 2023.
- Yang, Y., Song, Q., Gao, Z., Wang, G., Li, S., and Zhang, X. StyDeco: Unsupervised style transfer with distilling priors and semantic decoupling. arXiv:2508.01215, 2025a.
- Yang, Z., Wang, H., and Hu, D. Efficient quantification of multimodal interaction at sample level. In *Forty-second International Conference on Machine Learning*, 2025b.
- Yazgan, A. M. The problem of the century: Brain rot. *OPUS Journal of Society Research*, 22(2):211–221, 2025.
- Yousef, A. M. F., Alshamy, A., Tlili, A., and Metwally, A. H. S. Demystifying the new dilemma of brain rot in the digital era: A review. *Brain Sciences*, 15(3):283, 2025.
- Yu, H., Jeong, S., Pawar, S., Shin, J., Jin, J., Myung, J., Oh, A., and Augenstein, I. Entangled in representations: Mechanistic investigation of cultural biases in large language models. arXiv:2508.08879, 2025.
- Zerkouk, M., Mihoubi, M., and Chikhaoui, B. A comprehensive review of AI-based intelligent tutoring systems: Applications and challenges. arXiv:2507.18882, 2025.
- Zhai, C., Wibowo, S., and Li, L. D. The effects of over-reliance on AI dialogue systems on students’ cognitive abilities: a systematic review. *Smart Learning Environments*, 11(1):28, 2024.
- Zhang, S., Yao, L., Sun, A., and Tay, Y. Deep learning based recommender system: A survey and new perspectives. *ACM computing surveys*, 52(1):1–38, 2019.
- Zhang, S., Chen, Z., Chen, L., and Wu, Y. CDST: Color disentangled style transfer for universal style reference customization. arXiv:2506.13770, 2025a.
- Zhang, Y., Sun, J., Feng, L., Yao, C., Fan, M., Zhang, L., Wang, Q., Geng, X., and Rui, Y. See widely, think wisely: Toward designing a generative multi-agent system to burst filter bubbles. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, pp. 1–24, 2024.
- Zhang, Y., Diddee, H., Holm, S., Liu, H., Liu, X., Samuel, V., Wang, B., and Ippolito, D. Noveltybench: Evaluating language models for humanlike diversity. In *Proceedings of the Conference on Language Modeling*, 2025b.
- Zhao, S., Hong, M., Liu, Y., Hazarika, D., and Lin, K. Do LLMs recognize your preferences? Evaluating personalized preference following in LLMs. In *Proceedings of the International Conference on Learning Representations*, 2025.
- Zhou, L.-Y. and Wang, Y.-Y. Simulation of personalized english learning path recommendation system based on knowledge graph and deep reinforcement learning. *Sci Rep*, 15(1):34554, October 2025.
- Zhou, X., Zhu, H., Mathur, L., Zhang, R., Yu, H., Qi, Z., Morency, L.-P., Bisk, Y., Fried, D., Neubig, G., and Sap, M. SOTOPIA: Interactive evaluation for social intelligence in language agents. In *Proceedings of the International Conference on Learning Representations*, 2024.
- Zhou, Y., Wang, H., He, J., and Wang, H. Review-based explainable recommendations: A transparency perspective. *ACM Transactions on Recommender Systems*, 3(3):1–20, 2025.
- Zhou, Z. Empathy in education: A critical review. *International Journal for the Scholarship of Teaching and Learning*, 16, 11 2022. doi: 10.20429/ijstol.2022.160302.

1100 **A. Appendix**

1101 **A.1. Frequently Asked Questions (FAQ)**

1102 **Q: Why are rights challenging to define?**

1103 **A:** Rights are challenging to define as different groups hold mutually exclusive conceptions of what counts as a right. Some  
1104 examples of how different groups define rights differently include, for example, how the prohibition against torture is  
1105 widely recognized as a human right under international law, yet its practical interpretation varies. Singapore permits judicial  
1106 corporal punishment ([United Nations Human Rights Office of the High Commissioner, 2019](#)), which is condemned by the  
1107 United Nations as a form of torture ([United Nations Human Rights Office of the High Commissioner, 1984](#)). Similarly,  
1108 according to the Universal Declaration of Human Rights, men and women of full age, without any limitation due to race,  
1109 nationality or religion, have the right to marry and to found a family ([United Nations, 1948](#)). However, same-sex marriage is  
1110 legally recognized in Taiwan ([Amnesty International, 2019](#)) and Sweden ([Government Offices of Sweden, 2018](#)), but not in  
1111 China ([Reuters, 2016](#)).

1112  
1113 As Wilson observes, rights are “ideologically promiscuous”, invoked across diverse political and moral positions without a  
1114 shared underlying definition ([Wilson, 2006](#)). Rather than treating rights as fixed universal morals, Wilson argues that they  
1115 are better understood in terms of what they do.  
1116

1117 **Q: How do you define “rights”?**

1118 **A:** As rights are challenging to define, this paper adopts the vernacular of “rights” as a term that identifies capabilities that  
1119 humans risk losing to AI assistants, namely the “rights” to choose, to own, to work, and to learn. This paper does not seek to  
1120 explicitly define what rights are, but instead presents descriptive research questions that shift the focus toward what AI  
1121 systems should do within expectations of agency regarding everyday processes. This aligns slightly to Wilson’s argument  
1122 that rights are better understood in terms of what they do, and we use this direction to guide the operationalization of AI  
1123 principles that safeguard agency when humans choose, own, work, and learn.  
1124

1125 **Q: Are the “rights” in this paper exhaustive, and what other “rights” can readers consider?**

1126 **A:** In this paper, we have focused on the “rights” to choose, own, work, and learn due to their relevance, pervasiveness,  
1127 and urgency. Simultaneously, we recognize that they do not exhaust the ways AI systems affect human agency. There  
1128 exist combinatorial “rights” (e.g., the “Right” to Learn at Work, the “Right” to Choose for Learners) and other rights that  
1129 warrant attention. The “Right” to Safety aligns with [Jobin et al. \(2019\)](#)’s non-maleficence principle and captures concerns  
1130 around deepfakes, scams, and other harms arising from AI-augmented manipulation. The “Right” to Dignity ensures that  
1131 communities retain agency and dignity over how data or models represent them. The “Right” to Access/Participate ensures  
1132 access to AI-supported opportunities, goods, and services, and inclusion of diverse stakeholders in co-designing AI systems.  
1133 Each “right” merits a dedicated in-depth discussion, hence demanding separate future works.  
1134  
1135

1136 **Q: Are the research questions in this paper exhaustive?**

1137 **A:** The research questions are not exhaustive. Formulating good research questions is itself a difficult and underexplored  
1138 problem, especially in an emerging area like AI and human rights. Rather than claiming to present an exhaustive set of  
1139 questions, this paper attempts to identify what kinds of questions are worth asking at this stage. Our goal is to make explicit  
1140 the lines of inquiry that matter, not to close them off.  
1141  
1142

1143 **Q: How were the research questions in this paper chosen?**

1144 **A:** The questions in this paper reflect a particular way of thinking about expectations of human agency and our capabilities  
1145 in everyday processes. These are socially enacted and practically consequential, and are meant to be illustrative rather than  
1146 complete. We see this work as a step toward a more systematic framework that could, in future work, help researchers  
1147 generate and assess relevant questions in a principled way. When such a framework is more fully developed, it may provide  
1148 guidance not only on which questions to ask, but also on how to derive new questions by following the same underlying logic.  
1149  
1150

1151 **Q: Why are insincere students excluded from the “Right” to Learn with AI? Shouldn’t all students be entitled to the**  
1152 **“Right” to Learn?**

1153 **A:** The “Right” to Learn with AI is premised on the learner’s willingness to engage meaningfully in the process of learning.  
1154

Students who are not sincere in their intent to learn are unlikely to exercise this “right” responsibly, as their use of AI would not be directed toward understanding or growth. *Even in the absence of AI, such learners may find ways to circumvent learning (for example, by copying peers), and thus fall outside the scope of this framework.* Our focus is therefore on learners who demonstrate genuine intent and who can benefit from AI systems designed to support learning.

**Q: The paper calls for a wide audience to work on research questions to develop AI that benefits humans. Is it practically feasible?**

**A:** It may appear naive to think that researchers, developers, social scientists, humanities scholars, policymakers, and users of AI would want to work on these research questions, when real decisions are usually shaped by companies, organizations, and incentives like profit maximization. However, our hope with this paper is also that, as public awareness grows regarding the trade-offs between AI assistance and human agency<sup>14</sup>, user expectations may increasingly shape the design of AI systems.

When users express sufficient concern regarding human agency or select competitor AI systems that support human agency, companies may be compelled to offer features that do not directly maximize profits, motivating work on the research questions mentioned in this paper. Such features include improved transparency, explainability, or the deliberate introduction of organizational “friction” in AI-supported workflows within the company. These measures can function as safeguards that preserve meaningful human oversight, even when full automation would be technically feasible. For example, Red Hat is now one of the world’s leading providers of enterprise open-source software solutions. Although open sourcing software does not inherently maximize profits, sustained consumer preference has helped establish the transparency brought about by open sourcing as a competitive norm. Similar dynamics may arise for AI systems as users become more aware of the balance between human agency and AI assistance, motivating firms to design to support human agency even if they impose additional cost or organizational constraint.

**Q: What are the promising extensions or future works of this paper?**

This paper poses research questions; each research question can extend into one or multiple future works. Furthermore, we also recognise that questions of enforcement in AI governance are inherently difficult, particularly when decisions are shaped by powerful economic and organizational incentives. Therefore, following operationalization of AI principles, a natural next step for future research from the perspective of AI governance is to map the proposed “rights” onto existing human rights frameworks that already have legal and institutional enforcement mechanisms, shifting the discussion from what developers and researchers should do to what can be meaningfully enforced. Doing so would likely require involvement beyond individual companies, and include governments and broader policy tools such as regulation and taxation, rather than relying on voluntary compliance alone. This is beyond the scope of this paper.

**Q: Since the paper is well-researched and well-cited, it is clear that each of the risks stems from an already-existing research program: some researchers are pursuing this research. Why is this paper still necessary?**

**A:** We highlight that the way research is currently conducted is often misaligned with its intended real-world impact. Recent large-scale analyses (Duede et al., 2024) show that AI has rapidly diffused across disciplines, but remains loosely integrated, producing “semantic tension” between AI-driven and domain-specific research. For example, AI-driven research may emphasize benchmarks and predictive performance, while domain-specific research prioritizes explanation or interpretability. As a result, even when addressing the same problems, AI-driven research may not properly address risks to AI use; while these lines of work coexist in both domains, they may not meaningfully integrate, limiting downstream adoption and impact. This is evidenced in domain-specific papers for medicine (Hassan et al., 2024), law (Magesh et al., 2025), finance (Aldasoro et al., 2025), and education (Tian et al., 2025).

Our paper is necessary as it reorients how research is conducted, to balance AI assistance and human agency, to benefit humans. It highlights risks to human agency reflected by widely recognized AI principles and operationalizes these AI principles through research questions that welcome interdisciplinary collaboration. The research questions hence explicitly accounts for domain context, to achieve operationalization and meaningful adoption and impact.

<sup>14</sup>Users may increasingly feel *urgency* (Section A.3.3) to balance AI assistance with human agency.

## A.2. “Rights” Framing

In this paper, we use the term “rights” to *articulate* the capacities that users should retain in choosing, owning, working, and learning in an AI-mediated world. By framing these capacities using the term “rights”, we can map risks that undermine human agency to high-level AI principles, and operationalize these AI principles using research questions or technical solutions that directly alleviate these risks.

### A.2.1. LEGAL RIGHTS VS OUR “RIGHTS”

The authors of this paper recognize that the term “right” carries significant legal and political weight. *The use of the term “rights” in this paper does not depend on treating the “Rights” to Choose, Own, Work, and Learn as enforceable legal entitlements or as claims within a formal rights framework.* The history of international human rights demonstrates that even widely accepted rights can be contested, and the creation of new rights requires extensive normative justification and political consensus (U.S. Mission Geneva, 2017; Seah, 2022). We do not claim to define new human rights, nor do we seek to intervene in debates over their enforceability or legal interpretation. Instead, we borrow the vocabulary of rights for its clarity and ability to capture what everyday users of AI expect to preserve in their interactions with AI. Framing these expectations as rights enables us to articulate what is at stake, to identify specific risks to user agency, and to motivate research directions that reinforce these core capacities. The phrasing thus serves a practical function: it provides a concise and accessible way to communicate the normative commitments underlying our technical proposals, without invoking the full legal and political apparatus of human rights discourse.

### A.2.2. FORCE OF NORMATIVITY

The “rights” to choose, own, work, and learn do not draw on the term “rights” for their force of normativity. Instead, they draw on well-established normative foundations, particularly the global synthesis of AI ethics principles surveyed by Jobin, Ienca, and Vayena (Jobin et al., 2019). Jobin et al.’s analysis comprises 84 ethics guidelines documents from various organizations including private companies, academic and research institutions, governmental agencies, inter-governmental and supra-national organizations. It maps the global landscape of AI ethics and reveals strong convergence on a set of high-level principles, including transparency, justice and fairness, non-maleficence, responsibility, and privacy. Although the terminology in this position paper differs (in using the term “rights” rather than “principles”), the “rights” draw on the global consensus from Jobin et al.: that AI development efforts should integrate with the AI principles through substantive ethical analysis and adequate implementation strategies. These will empower users, protect their agency, and promote equitable, trustworthy, and beneficial outcomes.

### A.2.3. “RIGHTS” VS AI PRINCIPLES

A key insight from Jobin et al. (2019) is that although AI ethics guidelines converge on shared values, they typically offer limited operational guidance for how those values should be implemented in practice (Jobin et al., 2019). Principles such as transparency or autonomy are often stated abstractly, without specifying how an AI system should behave or what technical mechanisms would make these principles actionable. The contribution of this position paper is in addressing this gap by reorganising these principles into “rights” and risks that can be resolved by answering the research questions proposed. The output-oriented research questions proposed translate high-level normative desires or expectations from Jobin et al. into concrete behavioral targets for AI systems. Rather than engaging in legal interpretation or philosophical argumentation which lie outside the primary expertise of scholars engaged in developing AI systems, this paper leverages their technical authority to show how these AI principles can be expressed as system requirements, research challenges, and design objectives.

## A.3. Selecting the “Rights”

We focus on these four “rights” for three complementary reasons. First, their **relevance** (Section A.3.1): they map cleanly onto established AI-ethics scholarship, particularly Jobin et al.’s synthesis of global AI-ethics frameworks, which shows strong convergence around principles such as transparency, justice and fairness, non-maleficence, responsibility, and privacy. Second, their **pervasiveness** (Section A.3.2): these “rights” track everyday human capacities that AI systems are increasingly insinuating themselves into, shaping how people learn, work, own, and choose in routine and often unnoticed ways. Third, their **urgency** (Section A.3.3): the erosion of these “rights” is already producing tangible social harms, and failure to address them now risks deeper, more systemic consequences as AI adoption accelerates. Taken together, these considerations

motivate our choice of daily processes as impactful both in breadth and depth as they affect large segments of society bearing on fundamental aspects of human agency and well-being.

### A.3.1. RELEVANCE

The four “rights” map cleanly onto established AI-ethics scholarship, particularly Jobin et al.’s synthesis of global AI-ethics frameworks, which shows strong convergence around principles such as transparency, justice and fairness, non-maleficence, responsibility, and privacy.

Transparency is foundational to all four rights because individuals cannot exercise agency over choosing, owning, working, and learning if AI systems operate opaquely. For the “Right” to Choose, transparency allows users to recognize when choices are being shaped or constrained by algorithmic curation rather than freely made. For the “Right” to Own, it enables tracing data provenance and authorship, which is essential for attribution and credit. In the “Right” to Work, transparency clarifies how AI aids in tasks in the workplace. Finally, for the “Right” to Learn, explainability helps learners understand answers and the reasoning process, preserving inquiry and reflection.

Justice and fairness ensure that AI does not systematically disadvantage particular individuals or groups across domains of everyday life. In choice, unfair recommender systems can systematically narrow exposure to viewpoints or opportunities for particular communities, undermining equal participation in cultural and civic life. In ownership, inequities arise when some creators experience appropriation of their work without recognition or compensation while others benefit. In work, unfair automation can disproportionately advantage certain roles or demographic groups. In learning, unfair AI tools can advantage some learners while penalizing others through unequal access to assistance or biased assistance.

Non-maleficence directly motivates the protection of these “rights” by foregrounding the harms or *risks* that arise when AI systems are deployed without restraint. Non-maleficence is most immediately salient in discussions of the “Right” to Choose, as this right directly confronts cases where AI systems are deliberately deployed by malicious actors to spread misinformation and disinformation, thereby causing harm to users’ decision-making autonomy. Non-maleficence is also relevant to the “Rights” to Learn and Work where overreliance can result directly from the over-use of AI systems for these tasks. Safeguarding these “rights” is thus a practical instantiation of the obligation to avoid foreseeable harm.

Responsibility and accountability are not only AI governance requirements but also consequences of human agency as to act meaningfully is to be answerable for actions and their effects. It is also most salient in discussions of the “Right” to Work where workers are often compelled to use AI tools and the resulting outcomes expose how responsibility is diffused across employers, developers, vendors, and regulators, frequently without clear delineation. It is also related the “Right” to Own where the irresponsible use of data, such as when training generative models to create art, affects the livelihood of human creators. Preserving these rights therefore requires maintaining clear links between agency, action, responsibility, and accountability.

Privacy and data protection are enabling conditions for human agency because the abilities to choose, own, work, and learn freely depend on having spaces that are not fully observable, recorded, or optimized against. When individuals know that their data, behavior, or creative outputs are continuously collected and analyzed, they may self-censor, conform, or avoid risk, hence diminishing their capacity to act autonomously. This is most obvious in the “Right” to Choose where pervasive data harvesting enables hyper-personalized persuasion that constrains genuine freedom of selection. Protecting privacy therefore sustains the agency required to choose, own, work, and learn freely.

For brevity, the mappings of the various principles to the rights are summarized in Table 1.

*Table 1.* The mappings of the various widely-recognized AI principles (Jobin et al., 2019) to the “Right” to Choose, Own, Work, and Learn. The ◦ indicates general relevance, while the ✓ indicates direct relevance, detailed in Section A.3.1.

AI Principles	Choose	Own	Work	Learn
Transparency	✓	✓	✓	✓
Justice and Fairness	✓	✓	✓	✓
Non-maleficence	✓	◦	✓	✓
Responsibility	◦	✓	✓	◦
Privacy	✓	◦	◦	◦

1320 A.3.2. Pervasiveness

1321 The four “rights” are pervasive because they correspond to capacities that are exercised continuously across the course of  
1322 ordinary life. It is difficult to imagine living without learning or working, whether through formal education and employment  
1323 or through everyday problem-solving and contribution. Likewise, creative expression is a natural human impulse (Boden,  
1324 2004) and the expectation that one can retain control or recognition over what one creates follows naturally. Lastly, choice is  
1325 pervasive through engagement with online platforms that curate information, shape attention, and influence preferences.  
1326 Together, these “rights” cut across age, occupation, and culture, making them not exceptional or aspirational ideals, but  
1327 foundational dimensions of everyday human agency.  
1328

1329 A.3.3. Urgency

1331 Across choosing, owning, working, and learning, the erosion of these “rights” is already manifesting in concrete and  
1332 observable ways, underscoring the urgency of addressing them before they become further entrenched.

1333 In moments of leisure, too, we increasingly surrender choice to AI algorithms that curate our For You Page, drawing upon  
1334 our browsing history, prevailing trends, and even contextual factors such as geographical location (Bhandari & Bimo, 2022).  
1335 In doing so, we place in AI’s hands a quiet but profound power: the authority to choose and shape what we see. This control  
1336 makes us vulnerable if there is little oversight or accountability to ensure that the content recommended is safe, constructive,  
1337 or free from harmful influence (Tommasele & Menczer, 2022; Bagchi et al., 2024). The danger deepens in an age when it is  
1338 often impossible to distinguish between AI-generated and authentic material (Gillham, 2024; Frank et al., 2024). Worse  
1339 still, AI-generated content itself is not neutral, reinscribing narrow stereotypes (Luccioni et al., 2023) and forcing audiences  
1340 to consume the same narratives while excluding diverse perspectives (Nicoletti & Bass, 2023). Human content is already  
1341 fraught with partiality, often designed to provoke reaction; when such dynamics are reinforced by algorithmic curation, the  
1342 risk arises of a vicious cycle where polarizing content is rewarded with attention and in turn amplified to attract ever more  
1343 views (Gillani et al., 2018). Preventing this spiral requires vigilance and transparency but will enable the preservation of  
1344 **The “Right” to Choose** the content we consume.  
1345

1346 Generative AI has also demonstrated a remarkable capacity to compose text, images, and music that both rival and enrich  
1347 traditional forms of art, expanding the imaginative powers of human creativity while lowering the barriers to artistic  
1348 expression. It claims to democratize art by increasing the reproducibility of art, transforming viewers from passive observers  
1349 into active participants or creators (Park, 2025). However, this very power sharpens the question of ownership: who rightfully  
1350 claims authorship of such works? Should it be the individual who crafts the prompt, the artists whose creations supplied  
1351 the training data, the subjects whose likenesses are reproduced, or the AI system itself as a novel creative agent? Without  
1352 a principled assignment of ownership, the problem quickly unfurls. Original artists may be disincentivized to produce  
1353 new works if their efforts serve only as fodder for AI training (Corrall, 2023); the commercial adoption of AI-generated  
1354 art also deprives artists of salaried employment, displacing them from the industry (Merchant, 2023); others may retreat  
1355 into privatization, restricting access to their art in fear of appropriation. These cases contradict the democratization of art  
1356 and place art further out of reach of ordinary people who are denied the opportunity to experience and appreciate artistic  
1357 expression in an open and shared domain. Safeguarding creativity in the age of generative AI therefore requires clear  
1358 recognition and protection of **The “Right” to Own**.  
1359

1360 In the workplace, AI is displacing fresh graduates (Simon, 2025) by automating white-collar entry-level jobs (Leopold,  
1361 2025). Office work once grounded in tacit know-how such as writing emails, analyzing data, and creating presentations  
1362 are now readily learned by models trained on vast corpora. Consequently, AI can elevate the performance of new and  
1363 lower-skilled employees even as it compresses the productivity premium of more experienced staff (Cazzaniga et al., 2024).  
1364 A caveat is that these gains are uneven: where baseline competence is already high, marginal improvements diminish  
1365 and skill decay accelerates (Macnamara et al., 2024; Brynjolfsson et al., 2023). Furthermore, widespread reliance on the  
1366 same generative tools tends to homogenize ideation, causing the new idea to converge if everyone draws from the same  
1367 engine (Anderson et al., 2024; Wadinambiarachchi et al., 2024). The net effect is a flattening of differentiation as the  
1368 genuinely productive work of junior employees is discounted as commonplace, while the distinctive contributions of highly  
1369 skilled workers are blunted by convergence toward median solutions. With the use of AI in the workplace, the distinctions  
1370 between individual and AI-generated contributions are blurred and the availability of work and jobs are limited, threatening  
1371 **The “Right” to Work**.

1372 It is easy, even facile, to become over-reliant on an ever-present and patient tutor, raising concerns regarding the dulling  
1373 of habits of inquiry and an erosion of learners’ cognitive abilities (Zhai et al., 2024); and because they require minimal  
1374

1375 effort (Bailey & Warner, 2024; Khan, 2024b), they become ripe for exploitation when completing assignments and  
1376 examinations (Goodier, 2025; Chan & Lo, 2025). learners, too, have begun to sound like their tutors (Yakura et al., 2025)–  
1377 systems trained on human language have incorporated their own cultural inflections, and the widespread use of such tutors  
1378 narrows the space for genuine self-expression. Left unchecked, this drift erodes linguistic and cultural diversity (Yakura  
1379 et al., 2025) and diminishes the articulation and formulation of independent opinions necessary for learning; and when  
1380 models are optimized for being correct, it tempts learners to trade the slow work of reasoning and the virtue of making  
1381 instructive mistakes for a shortcut to the correct answer. These put into jeopardy **The “Right” to Learn**.

#### 1383 **A.4. The Need for Operationalization**

1384 Beyond addressing existing gaps in the literature on operationalizing AI governance frameworks, we also offer an additional  
1385 motivation for operationalization: that technology is not value-neutral. Values are already embedded in the systems we build  
1386 and use, and making them explicit allows us to shape their effects rather than leaving them implicit and unexamined.

1388 Technologies are often treated as neutral tools whose social consequences arise only through use. However, a long  
1389 line of work argues that technologies are never fully value-neutral: even when designers intend neutrality, values enter  
1390 through design choices, assumptions, and the social arrangements required for deployment. As argued by Winner, certain  
1391 technologies such as nuclear power are inseparable from specific forms of social organization, authority, and hierarchy that  
1392 must exist for them to function at all (Winner, 1980). In this sense, technologies cannot be understood as apolitical artifacts,  
1393 but as systems that embody and enact values. This observation motivates our approach: rather than debating which values  
1394 technologies should reflect, we argue that embedded values must be made explicit and operationalized. Doing so is necessary  
1395 not merely because technologies have social impact, but because values are already present—whether acknowledged or  
1396 not—and therefore should be systematically examined and designed for.

## B. Illustrative Solutions

### B.1. The “Right” to Choose

#### B.1.1. EXAMPLE FORMULATION OF AN OPTION SET

We define the set of options that AI systems can recommend to users as an *Option Set*.

**Definition B.1.** Let the *Option Set* be a finite set

$$\mathcal{O} = \{o_i\}_{i=1}^n,$$

where each option  $o_i$  represents a possible source of information relevant to a decision problem,  $Y$  (for example, a task variable, a social media reel for a user with preferences, a result on a google search). The cardinality or number or *quantity* of available options is  $|\mathcal{O}|$ . A valuation function

$$\text{val}_Y : \mathcal{O} \rightarrow \mathbb{R}_{\geq 0},$$

maps each  $(o_i)$  a non-negative value  $\text{val}_Y(o_i)$  that represents its informativeness with respect to task  $Y$ .

The *quality* of the options should consider how informative  $\mathcal{O}$  is to  $Y$ .

While maximizing

$$\sum_{i=1}^n \text{val}_Y(o_i)$$

may appear to yield an optimal option set since it aggregates the most informative options for the task, such an approach assumes that users will process the entire set of options presented to them. In practice, users may attend to only a subset of options, select a single alternative, or experience cognitive fatigue when confronted with overly dense or information-heavy sets. Moreover, evaluating options solely based on their individual informativeness overlooks interactions among options. As a result, it is not straightforward to curate an optimal option set, and highlights the need for more deliberate option-set curation.

For example, AI solutions that craft the option set can examine the relational structure among all or the top-k options, such as how distinct or complementary they are, and how they collectively frame the decision space. Understanding these relationships could reveal how to construct option sets that maximize informativeness and preserve diversity and interpretability, thereby improving decision making.

#### B.1.2. ILLUSTRATIVE SOLUTIONS FOR RQ1: HOW CAN AI SYSTEMS PRESENT DIVERSE AND CREDIBLE CONTENT THAT EXPANDS USERS’ INFORMATIONAL EXPOSURE, WHILE PRESERVING MEANINGFUL AND SUSTAINED ENGAGEMENT?

AI systems structure and present limited sets of options for users to choose from (Mills & Sætra, 2024; Agan et al., 2023), and can be adapted to present diverse and credible content. For example, they can surface left-, center-, and right-leaning takes on the same story (AllSides, 2025) to encourage diversity or can integrate third-party fact-check labels to encourage credibility (Papakyriakopoulos & Goodman, 2022).

##### On presenting diverse content

AI systems may present diverse content to users *passively* by presenting alternative viewpoints to users, but not requiring users to explore them. For example, information may be clustered into multiple perspectives (AllSides, 2025), and users can select the viewpoints they are interested in, knowing that alternative viewpoints are available. Alternatively, AI systems may present diverse content to users *actively* by deliberately guiding users through a range of perspectives before being presenting information most aligned with their preferences. For example, a conversational agent could introduce users to contrasting viewpoints before summarizing content most relevant to their search history (Zhang et al., 2024).

A promising direction to better address the diversity aspect of this research question is to develop quantitative metrics for measuring content diversity and integrating them into AI systems. One such metric could utilize information theoretic frameworks such as Partial Information Decomposition (Williams & Beer, 2010; Griffith & Koch, 2014; Bertschinger et al., 2014; Ince, 2017) to quantify unique contributions from different sources of information. This has also been extended to multimodal domains (Yang et al., 2025b; Liang et al., 2023). Another is Pyrorank (Kilitcioglu et al., 2023), which increases

1485 user-based diversity and mitigates systemic bias that traditional recommender system models learn from personalized data.  
1486 Ultimately, we propose a direction where diversity is mechanistic and controllable to expand a user’s exposure to different  
1487 sources of information.

1488 **On presenting credible content**

1489 While presenting diverse content, AI systems must also ensure that this content is credible. This domain of discerning  
1490 facts and biases has been a rich and active domain of AI research. For example, third-party fact-check labels can flag  
1491 disputed or misleading claims to create a direct barrier against the spread of false or biased content (Papakyriakopoulos &  
1492 Goodman, 2022). Counterfactual ranking tests, which estimate how different ranking policies would perform if content  
1493 exposure were changed, can also be used. These tests help identify and correct for hidden feedback loops or popularity-based  
1494 amplification in recommendations made to a user that might privilege extreme or biased views (Buchholz et al., 2024).  
1495 Finally, aggregating independent reports (such as crowdsourced fact-checks or external credibility signals) and weighting  
1496 them by the reputations of their sources aids in separating verified factual information from opinion or propaganda (Solovev  
1497 & Pröllochs, 2025). Together, these methods have made significant progress to reduce biases within AI systems in content  
1498 selection and presentation, fostering a more balanced informational environment. However, biases can evolve over time  
1499 and emerge from multiple factors—such as demographic features—interacting with each other (Shah & Sureja, 2025). To  
1500 reliably select and present credible content, developing frameworks or adaptive objectives to handle these emergent and  
1501 dynamic biases is a promising direction for AI research.

1503 **B.1.3. ILLUSTRATIVE SOLUTIONS FOR RQ2: HOW CAN WE DESIGN AI SYSTEMS THAT USERS CAN TRUST?**

1505 The black box nature of AI systems prevents humans from fully understanding and interpreting its workings, even for  
1506 computer scientists or those with specialized training (Burrell, 2016). Building meaningful reasoning and justification and  
1507 communicating the intent behind AI outputs recommended to users can help them build trust in AI systems. Prior works have  
1508 partially addressed the transparency concerns in this research question by including explainable recommendation models that  
1509 provide differentiated levels of model transparency to explain how individual user preferences shape recommendations (Zhou  
1510 et al., 2025) and mechanisms to seek simple (low complexity) and effective (high strength) explanations for the model  
1511 decisions (Tan et al., 2021). However, current works have not sufficiently address the “trust” aspect of the system to make  
1512 more informed choices. Surveys have shown that both under- and over-transparency can erode trust in the AI system’s  
1513 recommendations, especially when users have previously disagreed with the AI system’s outputs (Kizilcec, 2016). A  
1514 possible solution is to introduce customizable and adaptive transparency mechanisms and safeguards that account for user  
1515 needs or the context of the situation.

1516 For example, older social media users face risks of misinterpretation or misinformation due to differences in demographics,  
1517 education, media literacy, or worldview (Brashier & Schacter, 2020), increasing their vulnerability to scams, health  
1518 misinformation, or ideological manipulation. In this case, custom system safeguards may include filtering and warning  
1519 mechanisms when emotionally manipulative content is detected. In these contexts, greater transparency may be necessary  
1520 to help users understand how their behaviors interact with these safeguards and enable them to engage with social media  
1521 without concern that the information presented is misrepresented.

1523 As another example, adolescents and young adults excessively exposed to low-quality online materials (e.g., Italian brain  
1524 rot) may experience emotional desensitization, cognitive overload, and a negative self-concept (Yousef et al., 2025). In such  
1525 cases, system-level safeguards that limit compulsive or excessive scrolling behaviors (e.g. doomscrolling) may be necessary,  
1526 alongside increased transparency to help users understand why certain interaction constraints or usage limits are applied, so  
1527 that they may engage with social media without being subject to manipulative interaction patterns that exploit attention or  
1528 emotional vulnerability.

1529 As a last example, social media users may become trapped in echo chambers or filter bubbles as recommender systems  
1530 increasingly adapt content feeds based on accumulated personal data and past interactions. In such cases, system safeguards  
1531 may allow users to reset their personal data and escape echo chambers through a transparent system that explicitly explains  
1532 the consequence of the user’s action (Tan et al., 2023). For example, before a user selects content, the system could explain  
1533 the potential consequences of this action on the user’s future content feeds. If these changes are undesirable, users may  
1534 choose not to proceed or provide explicit guidance to keep their content-feed unchanged, enabling users to regain control  
1535 over their content feeds and removing undesirable recommendations.

1540 **B.1.4. ILLUSTRATIVE SOLUTIONS FOR RQ3: HOW CAN AI SYSTEMS EMPOWER USERS TO REJECT, ADJUST, OR**  
1541 **WITHDRAW FROM ALGORITHMIC CURATION?**

1542 AI systems can empower users to directly adjust algorithmic curation. For example, obfuscation techniques can be used to  
1543 block inference of user preferences (Slokom et al., 2021). Such an AI system that supports user-led “re-personalization” or  
1544 personalization “resets” can, for example, help users generate videos that they may watch to reset a platform algorithm  
1545 curated to them.  
1546

1547 Existing approaches also adjust algorithmic curation by predicting user preference shifts over time. Preference shifts can be  
1548 predicted by modeling preferences as sequences (Zhang et al., 2019), or by trying to understand why preferences change  
1549 and how such changes influence user interactions with algorithmically curated systems (Wang et al., 2023). While these  
1550 works have contributed greatly towards curated and highly personalized algorithms, their objectives are often aligned to  
1551 retain user interactions and improve the user’s experience.  
1552

1553 As such, future works should consider addressing, more intentionally, the “user control” aspect of this research question  
1554 by incorporating insights from multiple disciplines. For example, researchers should incorporate insights from behavioral  
1555 researchers who study users’ mental models of recommender systems (Ngo et al., 2020) to inform the design of user  
1556 interfaces. This would allow users to better understand how recommendations are made and enable them to curate their  
1557 algorithms. Building on this, future works could integrate these user signals as a source of feedback to tune AI systems.  
1558 For example, AI systems can inform users when their feed becomes increasingly personalized, and allow users to signal or  
1559 articulate the desire to reduce personalization (Wang et al., 2022).

1560 AI systems can also support users to withdraw from algorithmic curation by providing neutral mode recommender views that  
1561 temporarily disable personalization and revert to population-level baselines for recommended content. This is in line with the  
1562 right for individual users to be forgotten (European Parliament and Council of the European Union, 2016). Recommender  
1563 views that temporarily disable personalization signal to users that they have withdrawn from algorithmic curation, and make  
1564 withdrawal a reversible, low-friction operation.  
1565

1566 **B.2. The “Right” to Own**

1567 **B.2.1. ILLUSTRATIVE SOLUTIONS FOR RQ4: HOW CAN THE AUTHORS OF AN AI-MEDIATED ARTWORK BE**  
1568 **IDENTIFIED AND THEIR CONTRIBUTIONS QUANTIFIED TO ASSIGN CREDIT AND OWNERSHIP?**

1570 AI-mediated artworks are created through a collaborative process between human artists and AI systems. Sustained  
1571 interactions among the artist, the AI system, and data, such as dataset curation to refine the model and iterative prompting,  
1572 entangle authorship and complicates credit and ownership allocation (Bomba & De Angeli, 2025).  
1573

1574 Beyond artists who are actively involved in the creative process, artists whose works are used as training data for AI tools  
1575 may also be considered contributors. These contributions to the generated output can be quantified through data attribution  
1576 methods, based on influence functions (Koh & Liang, 2017; Mlodozieniec et al., 2025), watermarking (Lu et al., 2025c), or  
1577 retrieval from a separate datastore (Min et al., 2024).

1578 However, many of these methods require open and known training data, which is often unavailable for large, proprietary  
1579 foundation models. Without access to the training dataset, identifying contributors becomes infeasible. This limitation  
1580 underscores the need for new legal frameworks to regulate the use, disclosure, and access of training data, to support clearer  
1581 attribution for AI-mediated artwork (Hacker, 2021).  
1582

1583 Naïvely allocating credit based solely on quantified contributions can overlook crucial nuances in creative work, especially  
1584 on the significance of contributions. For example, contributions that lead to unique artistic expressions can carry more  
1585 weight, even if the measurable changes are small, and should be allocated more credit accordingly. Development of metrics  
1586 to account for these qualitative factors requires collaboration with artists, to leverage their expertise and judgment to ensure  
1587 that credit allocation reflects both process and creative value. Structured credit allocation frameworks can also provide  
1588 transparent and standardized ways to recognize contributions, via more finegrained roles and credit assignment, inspired  
1589 by contributor role taxonomy for research (Brand et al., 2015). However, some contributions arise from interactions so  
1590 interdependent that credit cannot be disentangled, and collaboration with game theorists could help model and fairly allocate  
1591 these collective credit.

1592 Credit allocation can inform ownership assignment, but the ownership of AI-mediated artworks remains legally ambiguous.  
1593 Engaging with legal experts is critical where revisions to legal standards that clarify ownership and provide tailored  
1594

protections for AI-mediated artwork can be explored (Watikinnakorn et al., 2023).

**B.2.2. ILLUSTRATIVE SOLUTIONS FOR RQ5.1: HOW CAN WE DETERMINE WHEN AI-MEDIATED STYLISTIC INFLUENCE DESERVES CREATIVE CREDIT?**

Recognizing and attributing intangible contributions in AI-mediated artworks may require assigning credit for AI-mediated stylistic influence. When such influence is acknowledged, credit may be attributed accordingly; when it is not, attribution can instead revert to the upstream or original creator. Attribution is conventionally credited based on expression, instead of style; therefore the question of whether stylistic influence merits credit should be informed by the perspectives of artists and legal experts. However, this determination can be facilitated by technical approaches that make style and expression more tangible and assessable:

**Style.** The concept of style in the research community has often been defined as surface-level characteristics like textures, color palette, and lighting, that can be disentangled from the semantic or compositional elements. Techniques such as Neural Style Transfer (Gatys et al., 2015), StyleGAN (Karras et al., 2019), and more recent methods such as SCFlow (Ma et al., 2025) exemplify how style can be disentangled. This aids artists and legal experts in clarifying whether, and in what form, each aspect of this disentangled style warrants attribution and credit.

**Expression.** In contrast, expression encompasses the overall creative realization of an artwork, including core components like semantic and compositional structures to express the creator’s intent. Techniques like scene-graph representations (Li et al., 2024b; Wang et al., 2024b; Dutta et al., 2025), vision–language alignment models (e.g., CLIP, GLIP) (Radford et al., 2021; Li et al., 2022; Liu et al., 2024d), and semantic disentanglement frameworks (Shuai et al., 2024; Dalva et al., 2024; Yang et al., 2023) illustrate ongoing efforts to model and quantify the semantic and compositional structures that form a key component of artistic expression. The next step is to seek input from artists and legal experts to clarify whether each AI-mediated artwork exhibits expression that warrants attribution and credit.

We provide an additional concrete machine learning formulation by demonstrating how creative credit can be awarded to AI using Shapley value computations. In the example, the AI credit is the Shapley value of the AI contributor in a two-player setup.

We consider a sentence which has been written by author A (such as a human): “The acquisition function selects promising hyperparameters for evaluation.” This sentence is then modified (with modifications in italics) by author B (such as an AI): “The clever acquisition function selects expensive, promising hyperparameters for evaluation.”. The Shapley values ( $\phi$ ) for authors A and B are, respectively:

$$\phi_A(v) = \frac{1}{2}(v(A) - v(\emptyset)) + \frac{1}{2}(v(A, B) - v(B)) = 0.5 \times 0.20 + 0.5 \times (0.75 - 0.65) = 0.15 \quad (1)$$

$$\phi_B(v) = \frac{1}{2}(v(B) - v(\emptyset)) + \frac{1}{2}(v(B, A) - v(A)) = 0.5 \times 0.65 + 0.5 \times (0.75 - 0.20) = 0.6 \quad (2)$$

Where  $v$  is the value function, and as a baseline, we set the empty sentence to have a value of 0. Therefore,  $v(\emptyset) = 0$ . The normalized human credit is  $\frac{0.6}{0.6+0.15} = 0.8$ . The normalized AI credit is  $\frac{0.15}{0.6+0.15} = 0.2$ .

In the experimental setup for this machine learning formulation, the value function  $v$  uses LLM-as-a-judge. The prompt, used on `gwen2.5:7b-instruct`:

You are an objective evaluator.

Your task is to assign a VALUE score to a sentence, focusing on the contribution of a specified source.

**DEFINITION OF VALUE:**

The value is the extent to which the specified contribution improves the sentence in terms of:

1. Semantic clarity
2. Informational usefulness
3. Appropriateness

1650 SCORING RULES:

- 1651 – Score must be between 0 and 1.
- 1652 – 0 = no contribution or harmful contribution
- 1653 – 1 = highly valuable contribution
- 1654 – Evaluate ONLY the contribution from the specified
- 1655 source .
- 1656 – Ignore contributions from other sources .
- 1657 – Prefer lower scores if uncertain .

1659 OUTPUT FORMAT:

1660 Return ONLY a single number between 0 and 1.

1661 -----

1662 Sentence : { sentence }

1663

1664 As this is an example machine learning problem, we also non-exhaustively suggest future directions for research:

1665

- 1666 • The number of contributors can be increased to more than 2.
- 1667
- 1668 • Multimodal contributions can be considered, such as text, speech, and images.
- 1669
- 1670 • Shapley value computations are compute intensive, especially if the contribution is computed per word or pixel or
- 1671 token, and if the contribution is longer (e.g. essays, movies). Therefore, approximations can be considered.
- 1672
- 1673 • LLM-as-a-judge can be replaced with a better metric, possibly in consultation with interdisciplinary expert (e.g.,
- 1674 lawyers, animators).

1675

1675 **B.2.3. ILLUSTRATIVE SOLUTIONS FOR RQ5.2: HOW CAN WE QUANTIFY THE EXPRESSIVE OVERLAP BETWEEN AN**

1676 **AI-MEDIATED ARTWORK AND THE ORIGINAL?**

1677

1678 The extent of expressive overlap between two or more artworks can help generative AI rebuke, or help artists prove, claims of

1679 copyright infringement. Recent works (Karras et al., 2019; Ma et al., 2025; Yang et al., 2025a; Zhang et al., 2025a; Qi et al.,

1680 2025) have demonstrated the feasibility of separating the semantic meaning of an image from its stylistic representation.

1681 Similarly, scene graph-based methods (Li et al., 2024b; Wang et al., 2024b; Dutta et al., 2025) capture relationships between

1682 visual entities and their compositional structure, providing a pathway to trace the creative expression that emerges from

1683 spatial arrangements. These methods can be extended to identify copyright-relevant features such as composition, subject

1684 framing, or symbolic motifs, that reflect the artist’s individual creative choices rather than generic stylistic conventions.

1685

1686 **B.2.4. ILLUSTRATIVE SOLUTIONS FOR RQ6: WHAT CAN AUTHORS DO TO PROTECT THEIR WORKS PRIOR TO**

1687 **LEGAL JUDGMENT?**

1688

1689 Authors currently have to rely on legal interpretation and subjective judgment to prove that their work was copied and have

1690 no tangible technical mechanism to detect, or demonstrate copying. For example, Singaporean photographer Zhang Jingna

1691 sued Luxembourg artist Jeff Dieschburg who had painted her 2017 Harper’s Bazaar Vietnam photograph and submitted it to

1692 an art competition, winning a prize. However, the case initially ruled against her, and this was only overturned two years

1693 later in 2024 (de Leon, 2024).

1694

1695 One attempt to protect works prior to legal judgment is to have creators embed watermarks within the creative process.

1696 These will enable creators to detect and assert ownership even when works are altered or their origins obscured as the

1697 watermark should persist in these derivative works. However, black-box and beige-box attacks involving diffusion models,

1698 noise injection, and leveraging semantic priors can achieve near-perfect watermark removal (95.7%) with negligible impact

1699 on the residual image’s quality (Shamshad et al., 2025). Consequently, ensuring the robustness of watermarking-based

1700 methods remains an open research challenge.

1701

1702 Authors may also have had their works “scraped” (lawfully or unlawfully, with or without permission) (Chesterman, 2024),

1703 but lack a tangible technical mechanism to detect this. In the absence of legal judgment, creators may instead rely on

1704 technical deterrence mechanisms to protect their works by introducing perturbations or poisoning into their works to disrupt

models trained on them. This can significantly degrading performance (Shan et al., 2024; 2023; Liu et al., 2024e) and

discourages unauthorized scraping of protected works for training datasets. However, recent work (Foerster et al., 2025) shows that images can be depoisoned, and protecting works from “scraping” remains an open research challenge.

### B.3. The “Right” to Work

#### B.3.1. ILLUSTRATIVE SOLUTIONS FOR RQ7.1: HOW CAN AI SYSTEMS DETECT/LEARN OPTIMAL MOMENTS TO INTERVENE IN A TASK (VS. LEAVING IT TO THE HUMAN)?

We can look towards two scenarios: one where an AI system selectively intervenes in a human worker’s task, and another where an AI system selectively defers to a human decision maker. In both, we want to balance AI assistance and human autonomy to achieve better task performance while preventing over-reliance on AI systems.

In the first scenario, humans are the main drivers in a task and AI systems can intervene when the worker’s performance drops. AI systems can detect signs of impasse such as repeated errors, prolonged inactivity, or inconsistent reasoning patterns that signal the need for AI assistance, perhaps due to fatigue. These signals are often multimodal. For example, signals from electroencephalography (EEG) and eye activity detectors are used for driver fatigue detection (Sikander & Anwar, 2018). Beyond static thresholds for intervention, AI systems can learn to assist the worker when the expected value of the agent’s action exceeds that of the human across defined goals (McMahan et al., 2024).

However, these methods optimize the AI to complement a fixed offline model of a human. This disregards, for example, the ability of the human to learn, or even a tendency for the human to detrimentally cognitively offload tasks (i.e., effects of cognitive offloading (Risko & Gilbert, 2016; Grinschgl et al., 2021; Cavicchi et al., 2025)), changing the optimal moment to intervene. Our work argues that solving RQ7.1 should also include behavioural experts such as psychologists, who can better balance the trade-offs between AI as a helpful assistant and human agency, and mitigate detrimental effects of AI use in the workplace. Such advice can also inform AI researchers in building models that benefit the user, and future work can better address this research question by integrating findings and perspectives from other disciplines to better quantify and balance out the positive and negative effects of AI intervention in these systems.

To give a proof-of-concept on how RQ7.1 can be translated into research problems for the ML community that can also involve interdisciplinary experts, we demonstrate how AI systems can learn optimal moments to intervene in a task by detecting user uncertainty. This work is inspired by existing works on verbalised confidence (Groot & Valdenegro Toro, 2024; Li et al., 2025). Using the corbt/all-recipes dataset, we create two sets of instructions per recipe for 500 recipes using the following prompts and save them as confident/unconfident instructions:

```
def make_confident_prompt(recipe: str) -> str:
    return (
        f"Recipe: {recipe}\n\n"
        "You are a senior chef verbally instructing a junior chef "
        "on how to prepare the given recipe. "
        "Ensure your verbal directives are clear, and include sensory "
        "cues the junior should look out for."
    )

def make_imposter_prompt(recipe: str) -> str:
    return (
        f"Recipe: {recipe}\n\n"
        "You are a senior chef verbally instructing a junior chef "
        "on how to prepare the given recipe. "
        "Avoid committing to clear directions, give vague or unsure "
        "measurements, unclear sequences, and overuse filler words."
    )
```

Then we combine the instructions (500 confident + 500 unconfident) and shuffle them, and use LLM-as-a-judge to classify whether the chef is confident. We achieve accuracy of 88.6%. A domain expert like a psychologist should advise on, e.g., other cases when, or how the AI system should intervene, and how to intervene in a way that can improve the performance of the less confident chef in the long run.

The second scenario occurs when AI systems are used to automate decisions and uncertain decisions must be deferred to a human expert. This ability to defer decisions to a human expert is often critical in high stakes settings such as healthcare and cybersecurity where experts can consider factors inaccessible to the AI (Gao & Yin, 2025), such as a patient’s medical history, to make an informed diagnosis. Current methods for expert deferral include learning a rejector (Charusaie et al., 2022), and a routing policy (Gao & Yin, 2025). Decisions can also be routed to multiple human experts (Keswani et al., 2021), and with personalized policies that accounts for each worker’s unique skill-sets and preferences (Bhatt et al., 2025).

**B.3.2. ILLUSTRATIVE SOLUTIONS FOR RQ7.2: HOW CAN AI SYSTEMS INTERVENE IN WAYS THAT COMPLEMENT THEIR USERS’ EFFORTS AND EXPERTISE?**

To intervene in ways that complement their users’ effort and expertise, AI developers must work with psychologists to understand how AI tools can improve work quality rather than merely completing tasks. A central insight from work psychology is that effective human-AI collaboration is not about substituting human judgment, but about designing assistance so that the combined human-AI system can achieve outcomes that neither could reliably achieve alone. Improved performance can emerge from how responsibilities, representations, and decision-making authority are distributed between human and AI (Woods, 1986).

An illustrative solution would be to design an AI system that is able to understand the social context, to know how to intervene to complement their users’ effort and expertise. The AI system can integrate multimodal signals, such as speech tone, facial expression, gesture, text, gaze, and environmental context, to infer intent and adaptively select actions to assist users while considering interpersonal nuance. Social perception facilitates human-AI collaboration (Puig et al., 2021), and recent multimodal foundation models can already fuse such cues across modalities, offering a basis for “socially grounded” learning (Lee et al., 2024). However, social intelligence in a workplace context also requires conceptual input from social and behavioral sciences, such as theories of communication competence (Coffelt et al., 2019), personality and affect in collaboration (Sackett & Walmsley, 2014), and situated cognition. Combining these perspectives with simulation environments such as SOTOPIA (Zhou et al., 2024; Wang et al., 2024a) could allow agents to learn and evaluate social reasoning in realistic, multimodal, and goal-directed settings.

Beyond low-level assistance, interventions can also be planned via dynamic human-AI team workflows where the AI system acts as a manager agent to decompose complex tasks given by humans, allocate and coordinate subtasks based on the complementary strengths of humans and AI, and adapt to changing preferences (Masters et al., 2025). Such workflow-level human-AI collaboration can lead to more optimal outcomes than those achievable by humans or AI systems alone.

**B.3.3. ILLUSTRATIVE SOLUTIONS FOR RQ8: WHAT STRATEGIES CAN HELP ORGANIZATIONS ENSURE THAT AI IS USED APPROPRIATELY AND SPARINGLY SO THAT EMPLOYEES CONTINUE TO DEVELOP THEIR OWN SKILLS?**

Drawing inspiration from cybersecurity practices, organizations can develop AI-use governance frameworks that balance innovation with accountability. In such a framework, employees periodically complete small and well-defined tasks that they are expected to perform without AI assistance. These tasks can serve as competency checks, analogous to white-hat penetration tests in cybersecurity. Employees who consistently rely on AI for tasks that should reflect personal expertise or domain understanding can be detected through watermarking (Lau et al., 2024; Lu et al., 2025a) or provenance tracking (Nikolic et al., 2025). This employees can then be enrolled in targeted retraining modules that reinforce skill development and responsible AI use. In doing so, this creates a culture where white-hat behavior (ethical and transparent use of AI) is rewarded, grey-hat behavior (borderline or excessive dependence) is monitored, and black-hat behavior (misuse or concealment of AI use) is subject to corrective education.

**B.3.4. ILLUSTRATIVE SOLUTIONS FOR RQ9: HOW CAN AI SYSTEMS SUPPORT FAIR RESPONSIBILITY ALLOCATION IN HUMAN–AI COLLABORATIVE WORK?**

Matthias (2004) identified a responsibility gap arising from AI systems as traditional frameworks are unable to assign moral or legal responsibility to either human users or AI developers. Recent works have attempted to address this by reframing responsibility as being diffused across multiple AI and human agents (Bleher & Braun, 2022), and formalizing it with a structured causal responsibility attribution framework to allocate responsibility in human-AI collaboration (Qi et al., 2024).

However, fair responsibility allocation remains challenging as it is dynamic and it also evolves alongside technological, social and legal developments. For example, responsibility allocation should be adjusted when AI systems communicate

1815 their uncertainty (Groot & Valdenegro Toro, 2024; Lau et al., 2025): if users act on highly uncertain AI recommendations,  
1816 more responsibility may fall on the user; conversely, if the AI system’s uncertainty is poorly calibrated or inaccurate, the AI  
1817 system should bear greater responsibility.

1819 **B.3.5. ILLUSTRATIVE SOLUTIONS FOR RQ10: HOW CAN AI SYSTEMS ENABLE A JUNIOR EMPLOYEE TO REACH**  
1820 **THE COMPETENCY OF A SENIOR EMPLOYEE?**

1821 One illustrative solution is a human-in-the-loop mentoring system that can turn everyday engineering workflows into training  
1822 and teaching signals. For example, when a junior employee writes a code review, incident report, or project proposal, the AI  
1823 system can suggest revisions based on corpora of high-quality senior examples, preserving tone, clarity, and diplomacy.  
1824 Senior human reviewers still sign off and their corrections are used to further refine the model’s coaching ability. An AI  
1825 system that complements the guidance of senior employees, by monitoring activities of junior employees in a team, detecting  
1826 misalignments with task goals, and providing targeted interventions, can enhance teamwork and improve overall team  
1827 performance (Seo et al., 2025). AI-supported personalized and real-time feedback can accelerate the development of junior  
1828 staff toward senior-level expertise.

1830 While more coaching data is important for training a better AI system, these data are often sensitive and cannot be shared  
1831 freely across organizations or teams. To enable aggregation of expertise without centralizing sensitive data, federated  
1832 learning provides a promising approach. Each organization can train locally on its own data, such as codebase, review  
1833 comments, incident post-mortems, design docs, and communication logs. A global model based on the compiled experience  
1834 of many senior employees is then updated via federated learning, sharing parameter updates rather than raw data (Kuang  
1835 et al., 2024; Shu et al., 2025). This lets the system learn general patterns of good engineering practice drawn from many  
1836 senior engineers that can be used to guide junior engineers.

1838 **B.3.6. WHEN CAN AI USE BE PREFERABLE IN THE WORKPLACE?**

1839 AI use can be preferable in specific workplace contexts where it complements or outperforms human efforts and expertise  
1840 by reducing cognitive burden, mediating complexity, or expanding the space of possible solutions. In particular, AI systems  
1841 can productively intervene to support communication, creativity, and complex trade-off reasoning.

1843 **RQ13. How can AI systems ease communicative loads and reduce communication barriers in the workplace?**

1844 Communication is central to collaboration but is often complicated by tone, emotion, language, and context. As AI systems  
1845 become integrated into workplaces, they can mediate interactions between co-workers to improve clarity and reduce friction.  
1846 Socially and culturally aware AI systems can recognize tone and sentiment, support more empathetic exchanges, and reduce  
1847 language barriers in multicultural environments. This support may take the form of multilingual capabilities, particularly for  
1848 scarce languages, such as those supported by SEA-LION (Ng et al., 2025)—as well as multicultural awareness that adapts  
1849 communication to diverse norms and preferences. By intervening in these ways, AI systems can prevent communication  
1850 breakdowns and foster more cohesive and effective collaboration.

1851 **RQ14. How can AI systems support creative problem solving and idea generation?**

1852 AI systems can also be beneficial in creative tasks when they expand, rather than collapse, the space of ideas. Instead of  
1853 defaulting to generic or standardized solutions, AI can support workplace brainstorming by generating diverse creative  
1854 pathways, building upon a worker’s initial ideas, and co-creating novel directions. This is particularly important given  
1855 evidence that current AI systems risk homogenizing outputs and suppressing originality (Kosmyna et al., 2025; Doshi &  
1856 Hauser, 2024). To mitigate this risk, we envision AI systems acting as collaborative colleagues that offer alternatives rather  
1857 than black-box oracles with gold-standard answers. Achieving this requires models with strong capabilities for generating  
1858 diverse and coherent ideas—a challenge even for current language models (Zhang et al., 2025b). Benchmarks such as  
1859 NoveltyBench (Zhang et al., 2025b) and metrics like the Creativity Index (Lu et al., 2025b) provide starting points for  
1860 evaluating and developing AI systems that facilitate idea generation without converging creativity.

1862 **RQ15. How can AI systems learn to help workers handle complex trade-off problems?**

1863 Many workplace tasks involve complex trade-offs where multiple goals and constraints must be balanced at once. Such  
1864 problems are cognitively demanding for humans, particularly when constraints are numerous, conflicting, or dynamic,  
1865 and when they involve both technical requirements and human values and preferences such as fairness, satisfaction, or  
1866 well-being. These also have multiple conditions that are hard to keep track of and weigh simultaneously, especially when  
1867 priorities change over time. In workplaces, these constraints are not purely technical and can include human values, fairness,  
1868 and satisfaction. AI systems can be clearly advantageous in such tasks involving complex trade-offs where multiple goals  
1869

and constraints must be balanced simultaneously. AI systems can assist by systematically tracking constraints, modeling trade-offs, and identifying balanced solutions that account for both hard rules and soft factors. For example, hospital scheduling requires balancing patient care needs, staff workload, and resource availability—an optimization problem that AI can handle more reliably while reducing human cognitive burden.

#### B.4. The “Right” to Learn

##### B.4.1. ADDITIONAL RISKS FOR SECTION 5

**RISK: Loss of Individual Agency, Reasoning, and Understanding.** Recent studies show that AI use initially enhances learning and higher-order thinking, but prolonged use can weaken cognition and knowledge retention (Sung et al., 2017; Tian & Zheng, 2024; Wang & Fan, 2025). Learning involves reasoning and recall, such as considering multiple sources of information, forming own conclusions, and identifying mistakes. By offloading these to AI, learners lose opportunities to hone these skills and engage meaningfully in learning without assistance, negatively impacting decision making and encouraging procrastination and “laziness” (Ahmad et al., 2023; Kosmyna et al., 2025). Furthermore, the immediacy of AI-generated responses shortcuts learning by making it feel effortless, hence discouraging engagement, independent understanding, and reasoning. This is exacerbated by educational technologies that frame learning as “effortless” (Mentutor, 2024; Shabanov, 2025). As learners offload learning to AI assistants and lose ownership of individual and independent understanding and reasoning, they cease to exercise individual agency (Agarwal, 2023; Gerlich, 2025; Kosmyna et al., 2025).

**RISK: Homogenization in Learning.** AI learning systems tend to promote homogenized forms of knowledge and expression that are disproportionately aligned with Western norms (Naous et al., 2024; Yu et al., 2025), raising concerns about the perpetuation of knowledge and expression shaped around a single dominant norm. Firstly, this diminishes learners’ abilities to interpret knowledge from their unique perspectives, conforming ideas generated across learners (Agarwal et al., 2025; Moon et al., 2025). These homogenization effects persist even after learners stop using AI models (Liu et al., 2024b). As learners lose the capacity to explore their own perspectives, learning ceases to be personally meaningful and engaging. Secondly, homogenized instructions can deter learners due to a possible cultural and linguistic mismatch. This undermines the learner’s sense of belonging and adversely impacts communication, engagement, and comprehension, which are essential for effective knowledge transfer (Fink, 2023; Gay, 2018; Steele & Cohn-Vargas, 2013).

##### B.4.2. ILLUSTRATIVE SOLUTIONS FOR RQ11: HOW CAN AI SYSTEMS BE DESIGNED TO CLOSE OR EVEN BRIDGE THE GAP BETWEEN TEACHERS AND LEARNERS?

Pedagogical research shows that teachers who can understand their students’ learning progress, empathise with, and get a greater knowledge of the individual needs of their students, can support their learning better (Zhou, 2022; Meyers et al., 2019). However, for teachers to go through all learner’s questions and responses from an AI system would be too tedious. Thus, such a solution can tap on existing literature on teacher-facing learning analytics dashboards (Kim et al., 2024), where researchers and teachers can craft modular or meaningful summaries of learner’s questions, and efficiently view the scope of responses given to them (Ryu et al., 2024; Peper et al., 2024).

Alternatively, for teachers who want to align AI support with upcoming lessons, we propose a simple classifier that can help teachers classify questions based on their topic. We attempt to classify questions from the HuggingFaceH4/MATH-500 dataset using the `qwen2.5:7b` model with the following prompt, and achieve an accuracy of 0.4380 on the 500 samples in the dataset. The confusion matrix can be found in Table 2.

```
def build_prompt(problem: str) -> str:
    topics_str = "\n".join(f"- {topic}" for topic in TOPICS)

    return (
        "You are classifying math contest problems by topic."
        "Choose exactly one topic from this list: \n\n"
        f"{topics_str}\n\n"
        "Return only the topic name. Do not explain.\n\n"
        "Problem:"
        f"{problem}"
    )
```

”Topic :”  
)

Predicted	A	C&P	G	IA	NT	PA	PC	NaN
Gold								
A	77	2	8	4	3	12	15	3
C&P	1	28	7	0	2	0	0	0
G	0	0	40	0	0	0	0	1
IA	54	1	11	9	5	0	9	8
NT	7	0	0	0	37	12	1	5
PA	19	11	20	1	8	16	7	0
PC	8	0	24	0	0	0	12	12

Table 2. Confusion Matrix. The classes include Algebra (A), Counting & Probability (C&P), Geometry (G), Intermediate Algebra (IA), Number Theory (NT), Prealgebra (PA), Precalculus (PC), and NaN.

### B.4.3. SUPPLEMENTARY RQS FOR RQ12

#### RQ16. How can AI systems measure a learner’s depth of understanding and dynamically calibrate its support accordingly?

Adjusting the depth and complexity of an AI system’s support to match a learner’s knowledge level can foster meaningful understanding (Zerkouk et al., 2025). Consider an AI system that first measures the learner’s understanding through interactive questioning and analyzing patterns in the learner’s explanations and revisions, and then interprets these signals using pedagogical tools and assessment frameworks developed in collaboration with social scientists like educators and pedagogists. From the measure, the pedagogist can control the AI system to dynamically provide information at the right level for the learner’s understanding. This can be done through a “knob” to calibrate the depth and complexity of the AI system’s responses by, for example, cultivating learning through dialog and explanation. This allows the AI to support the access to and the understanding of knowledge by leveraging the pedagogist’s expertise, and to uphold the Right to Learn. See Section B.4.4.

#### RQ17. How can AI systems strategically allow learners to make mistakes and guide them in examining assumptions and identifying errors to develop robust reasoning?

Making and identifying mistakes allows learners to detect inconsistencies between their mental models and new information, triggering repair of inadequate mental models and improving understanding for learning (Chi, 2000). Consider an AI system that can strategically intervene or withhold correct answers during learning, so as to guide learners to ask clarifying questions, examine assumptions, and reflect on counterfactual scenarios. This helps learners think from alternative perspectives, prompting them to notice inconsistencies, gaps, ambiguities, and assumptions or mistakes in their reasoning. Pedagogical expertise can provide guidance on the timing and form of these interventions and when mistakes should be revealed. Psychologists can also advise how these interventions can be introduced. By engaging learners in these forms of inquiry and reflection and forcing them to take a more effortful path towards construction of knowledge, AI systems can aid reasoning and understanding and help learners take active responsibility and practise individual judgment when learning. See Section B.4.5.

#### RQ18. How should AI systems allow learners to choose their learning strategy and align their reasoning and explanations with this choice?

Learners benefit when they are given control over learning strategies or solution paths, provided that the system constrains, guides, or contextualizes those choices (Zerkouk et al., 2025). For example, in learning mathematics, learners may prefer specific mental models such as algebra, model-drawing (Ng, 2022), or intuitive reasoning. Pedagogical expertise is required to identify the range of diverse and valid strategies. By reasoning from the learner’s perspective, AI systems can adapt to the learner’s understanding and context rather than impose standardized explanations, keeping learners engaged and preserving individual agency. Learners from different cultural contexts may also want to choose different learning strategies. Consider an AI system that is designed to accept diverse cultural perspectives when developed alongside pedagogical experts from diverse cultural contexts. This can prevent learners from converging to the same way of thinking or from being shaped by a single dominant perspective embedded in the AI, preventing intellectual homogenization. See Section B.4.6.

**B.4.4. ILLUSTRATIVE SOLUTIONS FOR RQ16: HOW CAN AI SYSTEMS MEASURE A LEARNER’S DEPTH OF UNDERSTANDING AND DYNAMICALLY CALIBRATE ITS SUPPORT ACCORDINGLY?**

Personalizing support such as instruction and feedback from an AI system to each learner’s needs (e.g., based on their level of understanding) can improve student performance (Zerkouk et al., 2025). For example, adaptive assessments (Vie et al., 2017) where questions and tasks are actively selected according to the learner’s responses, can be leveraged to lead to more efficient and accurate assessments of a learner’s understanding. AI systems can enhance adaptive assessments via retrieval-augmented generation (RAG) (Lewis et al., 2020), by collating a comprehensive knowledge base with assessments of various difficulty levels. At each turn, generating questions that are appropriately challenging can gather more information on the learner’s understanding level. Rather than numeric metrics that often miss qualitative nuances, AI systems can possibly also capture the learner’s depth of understanding more holistically, via implicit representations in personalized models (Ning et al., 2025).

These approaches to calibrate an AI system’s support to a learner’s knowledge must be grounded in existing pedagogical frameworks. For example, scaffolding-guided training for pedagogical alignment has been shown to provide adaptive guidance to learners of different levels (Liu et al., 2025) and can be used by pedagogists to personalize the treatment of learners after evaluating their level of understanding. Pedagogists may also evaluate assessments following cognitive depth taxonomies (Anderson & Krathwohl, 2001; Black & Wiliam, 1998) to personalize effective learning objectives for learners. AI systems can then be aligned with the identified learning objectives via personalized alignment (Chen et al., 2024) to the pedagogist’s recommendations.

Another interesting problem for the research audience is that as a learner acquires new knowledge, the data distribution shifts and the AI system must continuously update its internal representations to match the learner’s level. This may interest researchers working on continual learning methods (Shi et al., 2025) where a model learns from new data without full retraining. Crucially, this adaptation should not be driven by the model alone. Pedagogists can play an active role in advising how the system ought to be calibrated to the learner at different stages of development. More interestingly, drawing on their understanding of learning trajectories, pedagogists may even anticipate how the data distribution is likely to shift, thereby informing when and how such updates should occur.

We note that learner modeling has been studied for many years and has already led to practical applications (Liu et al., 2024c; Gao et al., 2021). However, this remains an area of interest as there is a gap between technical operationalization and pedagogically grounded outcomes (Oh & Ahn, 2024), and existing works point to mixed results about the effectiveness of AI-based tutoring systems (Zerkouk et al., 2025).

AI systems must therefore partner with pedagogical insights to meaningfully calibrate to a learner’s depth of understanding. For example, AI systems can be used to collect quantitative metrics like test scores at scale from students. However, the AI, possibly trained on inputs from a pedagogist, should interpret the scores to discern the learner’s depth of understanding or evaluate other qualitative metrics on student performance based on their experience. These can then be used to dynamically calibrate the AI system or set priors regarding learners (like adjusting a “knob”).

Existing works, while done with good intention, leave pedagogical experts out of the loop (Oh & Ahn, 2024). For example, they may focus on improving the accuracy of learner modeling, but leave the connection between model representations and how educators define or intervene on learning under-specified (Liu et al., 2024c; Gao et al., 2021).

**B.4.5. ILLUSTRATIVE SOLUTIONS FOR RQ17: HOW CAN AI SYSTEMS STRATEGICALLY ALLOW LEARNERS TO MAKE MISTAKES AND GUIDE THEM IN EXAMINING ASSUMPTIONS AND IDENTIFYING ERRORS TO DEVELOP ROBUST REASONING?**

Pedagogical theory suggests that making and identifying mistakes allows learners to detect inconsistencies between their mental models and new information, triggering repair of inadequate mental models and improving understanding for learning (Chi, 2000). A possible solution is to align AI tutor systems to teachers instead of students. AI systems can align to teachers with the desired pedagogical abilities using reinforcement learning, which is widely used for LLM alignment (Ouyang et al., 2022; Rafailov et al., 2023; Lambert et al., 2024). Reward signals can be collected in collaboration with educators and pedagogists, where responses with contrastive dialogues that invite reflection of plausible alternatives and surface assumptions can be scored higher than responses that directly provide an answer. AI systems trained on these rewards will be more pedagogically aligned to the teacher.

A pedagogist may also be able to advise on intervention timing and intervention strategies to improve learning outcomes.

Teaching strategies like withholding from correcting a learner’s mistakes and guiding learners to reflect and identify their own mistakes are capabilities that pedagogists are familiar with. AI systems will have to be developed to engage in multi-turn dialogues and to plan conversations, selecting responses based on pedagogical expertise that are most likely to lead to the learning goal after multiple turns (Chen et al., 2025). Pedagogically-aligned evaluations, which include assessing if AI systems can delay answer reveal (Maurya et al., 2025), are also crucial to ensure that AI systems are exhibiting the correct teaching behaviors. These evaluations can be certified by pedagogists and periodically updated, to keep in line with the latest pedagogical frameworks.

Researchers from multimodal (visual, audio, etc) AI can also be engaged in this effort. Pedagogical works suggest that learners benefit from well-designed multimodal input (Paivio et al., 1968). Multimodal cues based on pedagogical expertise can also be used to guide learners towards reflective diagnosis of their own errors. For example, adding graphical hints to diagrams as feedback when learners incorrectly solve a problem can successfully guide learners towards the correct solution (Rouinfar et al., 2014; Wu et al., 2025). Another example is syntax highlighting (Klock & Jan, 1986) where code elements are displayed in distinct colors to indicate their syntactic roles. When a segment appears in an unexpected color, the user immediately perceives an anomaly and can infer that something is wrong. Similarly, when teaching text-based reasoning, the system can signal logical or structural inconsistencies by color-coding arguments, counter-arguments, or key claims, without revealing the correct interpretation (Ruiz-Dolz et al., 2025; Liu et al., 2023). As with the other imagined solutions above, this should be developed alongside pedagogists, for example by engaging them to identify which logical or structural inconsistencies to signal, or what aspects of a text should be highlighted to guide learning.

We note that guided educational support (Liu et al., 2024a), such as through game-ification of learning (Kanervisto et al., 2025), has been studied for many years and has already led to practical applications. However, recent works suggest that the role of teachers and pedagogists in developing such educational support remains limited (Tan et al., 2025). This has highlighted a gap in research addressing the AI development needs of teachers as they integrate AI technologies into their teaching practices and exploring how AI technologies can be applied in education from both the perspectives of student learning and teacher instruction (Tan et al., 2025).

**B.4.6. ILLUSTRATIVE SOLUTIONS FOR RQ18: HOW SHOULD AI SYSTEMS ALLOW LEARNERS TO CHOOSE THEIR LEARNING STRATEGY AND ALIGN THEIR REASONING AND EXPLANATIONS WITH THIS CHOICE?**

Learning strategies differ across contexts. For example, model-drawing is the method of choice for teaching mathematics in Singapore (Ng, 2022), but other educational systems may prefer specific mental models such as algebra or intuitive reasoning. Recent works suggest a gap in research addressing how AI technologies can be applied in education from both the perspectives of student learning and teacher instruction (Tan et al., 2025), and many works that operationalize AI tools for learning rely heavily on evaluating learning strategies or learner performance via model-estimated learner progress rather than externally validated learning gains with pedagogists or teachers (Zhou & Wang, 2025). As a result, current works risk ignoring pedagogical input regarding which learning strategies may be available or even optimal for the learner.

To overcome this, one illustrative solution is to explicitly ask learners how they would like to learn, to enable learners to customize their learning experience and be active participants in the learning process. AI systems should have a list of pedagogically validated learning strategies to rely on depending on the learner’s choice, and should remember and adhere to the stated preferences. As current AI systems may experience difficulty following stated user preferences in a zero-shot setting (Zhao et al., 2025), current methods to improve preference following include prompting strategies that remind AI systems to consider previously stated preferences (Zhao et al., 2025), and extracting relevant information from history using RAG (Lewis et al., 2020). Preferences, such as on which learning strategy to use, can evolve dynamically, and AI systems should capture and align to these evolving preferences, especially in long-context conversations (Jiang et al., 2025).

The choice of learning strategy may also be context or culture dependent. If the learner does not explicitly request a certain learning strategy, then given the profile information of a learner, implicit choices such as cultural preferences can also be inferred, and AI systems can adapt their learning strategies accordingly. AI systems that are culturally aware and encapsulate cultural differences can mitigate bias towards specific cultural perspectives by aligning their responses better with the learner’s cultural context (Li et al., 2024a). This would allow the AI system to choose the best, pedagogically validated, learning strategy for the learner.

Sometimes, the pedagogist or teacher may choose to customize the learning experience on behalf of the learner, based on their teaching experience and expertise. To avoid duplicating the preceding argument with only role labels altered (i.e., changing learner to teacher), we provide a simple proof of concept. We consider a teacher using a large language model

2090 to teach mathematics to grade 5 students. The teacher carefully writes a prompt detailing how the large language model  
 2091 should explain the answer of various mathematics questions to a grade 5 student, bearing in mind that the student only  
 2092 has knowledge of introductory algebraic thinking and is fluent in multi-digit multiplication and division. We compare the  
 2093 response to the teacher’s prompt against a basic/ baseline prompt. We observe that such teacher prompts are useful for  
 2094 structuring (i.e., curating the way that information is presented to the student) mathematical solutions. We then use an  
 2095 LLM-as-a-judge to see which response better aligns to the needs of a grade 5 student. This proof of concept was performed  
 2096 on the GSM8k dataset, and we found that the LLM-as-a-judge (flan-t5) preferred the response to the teacher’s prompt 78.7%  
 2097 of the time ( $n = 1284$ ).

2098 The baseline and teacher-aligned prompts are as follows:

```

2100     def make_baseline_prompt(question: str) -> str:
2101     return (
2102         f"Question: {question}\n\n"
2103         "Explain step-by-step for a grade 5 student. "
2104         "End your response with 'Final answer: <answer>'."
2105     )
2106
2107     def make_aligned_prompt(question: str) -> str:
2108     return (
2109         "You are a math teacher."
2110         "Explain step-by-step for a grade 5 student.\n\n"
2111         "To align with the learning strategy of a grade 5 student,"
2112         "Follow this EXACT structure:\n\n"
2113         "[Understanding]\n"
2114         "Restate the problem in simple terms.\n\n"
2115         "[Plan]\n"
2116         "Explain what strategy we should use.\n\n"
2117         "[Guided Thinking]\n"
2118         "Ask 1 short question to guide the student "
2119         "(DO NOT answer it yet).\n\n"
2120         "[Step-by-step Solution]\n"
2121         "Solve the problem step-by-step.\n\n"
2122         "[Final Answer]\n"
2123         "Write: Final answer: <answer>\n\n"
2124         f"Question: {question}"
2125     )
2126
    
```

#### B.4.7. PRELIMINARY PROGRESS BY CHATGPT’S STUDY MODE AND GEMINI’S GUIDED LEARNING IN ADDRESSING RQ16-RQ18

2130 The introduction of Gemini’s Guided Learning (Heymans, 2025) and ChatGPT’s study mode (OpenAI, 2025) represents  
 2131 a step in the right direction. For example, ChatGPT’s study mode adapts instruction to learners’ proficiency levels using  
 2132 diagnostic questions and conversational history, addressing RQ16. Likewise, Gemini’s Guided Learning promotes critical  
 2133 thinking through guided questioning, aligning with RQ17. These examples underscore the relevance of the proposed  
 2134 research questions.

2135  
2136  
2137  
2138  
2139  
2140  
2141  
2142  
2143  
2144