

# Undermining Mental Proof: How AI Can Make Cooperation Harder by Making Thinking Easier

Zachary Wojtowicz<sup>\*†</sup> and Simon DeDeo<sup>‡</sup>

January 7, 2025

## Abstract

Large language models and other highly capable AI systems ease the burdens of deciding what to say or do, but this very ease can undermine the effectiveness of our actions in social contexts. We explain this apparent tension by introducing the integrative theoretical concept of “mental proof,” which occurs when observable actions are used to certify unobservable mental facts. From hiring to dating, mental proofs enable people to credibly communicate values, intentions, states of knowledge, and other private features of their minds to one another in low-trust environments where honesty cannot be easily enforced. Drawing on results from economics, theoretical biology, and computer science, we describe the core theoretical mechanisms that enable people to effect mental proofs. An analysis of these mechanisms clarifies when and how artificial intelligence can make low-trust cooperation harder despite making thinking easier.

## 1 Introduction

The widespread availability of generative artificial intelligence means that anyone can now cheaply and convincingly simulate the output of human mental effort across an unprecedented variety of tasks. This promises numerous benefits across nearly every aspect of society, but it also has begun to disrupt an equally broad set of social practices, such as sincere apologies (Glikson and Asscher, 2023), college assessment (Fitria, 2023; Cardon et al., 2023), online dating (Wu and Kelly, 2020), and wedding vows (LaGorce, 2023).

In light of these developments, many have come to see the technology as a double-edged sword: artificial intelligence cuts the cost of thinking, but it also—and, as we will argue, *for that very reason*—threatens vital elements of the existing social fabric, such as trust (Glikson and Woolley, 2020), privacy (Jain et al., 2023), public safety

---

<sup>\*</sup>Corresponding Author.

<sup>†</sup>Department of Economics and Harvard Business School, Harvard University, Cambridge, MA, United States; [zwojtowicz@fas.harvard.edu](mailto:zwojtowicz@fas.harvard.edu)

<sup>‡</sup>Dept of Social and Decision Sciences, Carnegie Mellon University, Pittsburgh, PA & the Santa Fe Institute, Santa Fe, NM; [sdedeo@andrew.cmu.edu](mailto:sdedeo@andrew.cmu.edu)

(Leslie, 2019), and even democracy itself (Allen and Weyl, 2024; Jungherr, 2023). Indeed, a majority of survey respondents now feel “more concerned than excited” about the increased use of artificial intelligence in their daily lives (Tyson, 2023).

Despite the urgency of these concerns, the scientific community has struggled to articulate general principles that explain why lowering mental costs through artificial intelligence undermines such a wide variety of seeming unrelated social practices. We may “know it when we see it” in any particular case, but the lack of integrative scientific frameworks has made it difficult to formulate general solutions to existing problems or foresee future harms.

In this paper, we highlight the key role that “mental proof” plays in facilitating cooperation in low-trust environments. Mental proofs are observable actions taken to certify unobservable facts about the minds who perform them. As we describe fully below, people use two distinct mechanisms to substantiate mental proofs: signaling theory (an idea primarily studied in economics and biology) and proof of knowledge protocols (studied in computer science). To function, both mechanisms rely on implicit assumptions about the cost structure of organic mental activity—assumptions that are rapidly being disrupted by the proliferation of artificial intelligence in daily life.

An appreciation for the structure of mental proofs therefore helps elucidate the underlying logic of artificial intelligence’s various social consequences. In highlighting the importance of mental proof and its relationship to thinking machines, our paper contributes to a wider, cross-disciplinary attempt to proactively understand and address the technology’s various social consequences (*e.g.*, Solaiman et al., 2023; Weidinger et al., 2021; Mirsky and Lee, 2021)

We illustrate the practical importance of mental proof in two “worked examples” drawn from everyday social domains: sincere apology and subculture formation. In both cases, we discuss how low-cost simulations of intelligent behavior undermine the efficacy of mental proofs and the vital social benefits they provide.

As our discussion makes clear, mental proof is most valuable in situations where honesty cannot easily be enforced. This implies that the category of harms we discuss will disproportionately impact those who are not already embedded in high-trust networks and formal institutions, thereby reinforcing existing structural barriers to social mobility and economic development. Our analysis does, however, suggest effective strategies for mitigating these deleterious consequences, and we conclude with a sketch of the framework’s implications for policy, technology, and everyday life.

## 2 Mental Proof

Communication can greatly enhance coordination. Our species’ remarkable capacity for coordination relies, at its foundation, upon an ability to reliably externalize the nuances of our internal mental states—not just our beliefs, but also our intentions, values, preferences, skills, understandings, commitments, abilities, *etc.*—in ways that others will not only understand, but trust. As has been pointed out by both behavioral scientists (*e.g.*, Tomasello et al., 2005) and philosophers of mind (*e.g.*, Gilbert, 1990; Bratman, 1992), people’s ability to share and understand intentions, in particular, is essential to the formation of collaborative acts that range from going on a walk together to drafting a new constitution.

Sometimes this is easy: in many contexts, we can simply speak our minds and others will have good reason to believe us. Despite its many conveniences, however, “cheap talk” breaks down when people have incentives to lie (Farrell, 1987). In such contexts, mere assertions lose their credibility, and even those who try to tell the truth will be dismissed.

A variety of social institutions help enforce honesty and thereby preserve the benefits of its coordinating function, most notably reputation (Fehr, 2009), norms (Elster, 1989), and formal punishment (North, 1991). Unfortunately, however, these institutions are not always available—*e.g.*, before reputation is established, when claims are difficult to verify, or in places where cultural and legal institutions are weak.

In these “low-trust” contexts, interacting parties must not only state, but endeavor to *prove*, claims about their minds. One way people furnish such proof is by taking observable actions which (given certain assumptions about the capabilities and structure of human brains) provide strong grounds for a relevant claim about the mind. We refer to such behaviors as constituting **mental proof**. The validity of mental proofs are primarily underwritten by two separate mechanisms: **signaling theory** and **proof of knowledge protocols**. We review each, in turn.

## 3 Mechanism One: Signaling Theory

Signaling theory was introduced into economics by Spence (1973) to explain how seemingly self-defeating behaviors (*e.g.*, knowingly pursuing a degree in a field one is unlikely to use) can still benefit rational agents by signaling information about one’s preferences or abilities (*e.g.*, that one is smart enough to graduate college) to others in a way that cannot be faked. Signaling was introduced to biology around the same time by Zahavi (1975) to explain an analogous class of animal traits and behaviors: phenotypes that seem to reduce fitness, such as peacocks growing elaborate tails or

gazelles stotting.<sup>1</sup> As the theory points out, these acts credibly communicate information to potential mates or predators precisely because of their self-handicapping effect; voluntarily wasting resources can credibly signal a wealth of resources to begin with.

The central insight of signaling theory is that rational agents only take actions they expect to be beneficial on net. A behavior can constitute definitive proof that the actor expects its benefits (including revealing information to you, the potential observer) to outweigh its costs. This pinpoints why signaling can only be accomplished by behaviors that incur real net costs when faked. The apparent downside of signaling behaviors are precisely what establish their credibility: the cost structure ensures that deceptive types cannot send the signal with impunity. This self-policing logic is what enables signaling equilibria to extend credible communication to low-trust environments, where honesty can not be enforced (by, *e.g.*, reputation).

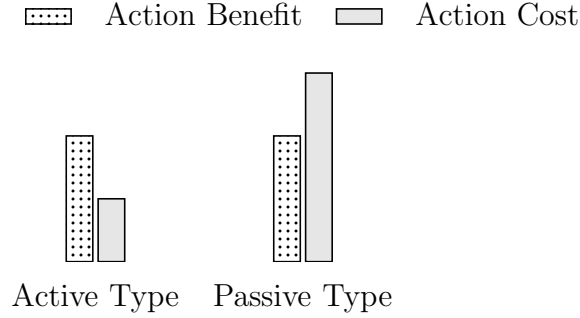
Figure 1 illustrates two focal cases of signaling equilibria. In the first case (subfigure a), two types of agents receive the same benefit from engaging in a signaling behavior, but incur different costs. From the perspective of an external observer, witnessing the behavior constitutes definitive proof that the agent is of the type with lower cost. Credibly communicating this distinction may benefit both parties, for example because having a low cost for the focal behavior (*e.g.*, finding mathematics easy) correlates with other traits (*e.g.*, abstract problem-solving ability) that are necessary for productive cooperation.

In the second case (subfigure b), both agents incur the same cost, but receive different benefits. From the perspective of an external observer, witnessing such a behavior constitutes proof that the agent is of the high-benefit type. Credibly communicating this distinction can profit both parties, for example because it reveals that the person will want the cooperative relationship to continue further into the future.

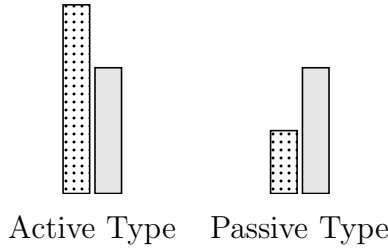
In economic domains, costly signaling has been used to explain a wide variety of phenomena, including lavish, uninformative advertising efforts (Milgrom and Roberts, 1986) corporate social responsibility initiatives (Cheng et al., 2014; Flammer, 2021), and disclosures made by ventures seeking early-stage funding (Ahlers et al., 2015). Signaling is not only confined to “western, education, industrial, rich, and democratic”

---

<sup>1</sup>“Stotting” refers to a behavior in which a quadruped acrobatically leaps into the air and contorts their body (Smith and Harper, 2003). Fitzgibbon and Fanshawe (1988) document, for example, that: Thomson’s gazelles are more likely to stott when wild dogs were nearby; that dogs were more likely to chase gazelles which stotted less frequently than their peers; and that gazelles which stotted more intensively had a better chance of escaping, conditional on being chased. This and other evidence has been taken to support the underlying hypothesis that such behaviors honestly signal a prey’s capacity to outrun predators, which makes them less likely to give chase. This cooperative détente between adversaries saves both energy (Caro, 1994).



(a) Equal Benefits  $\Rightarrow$  Different Costs



(b) Equal Costs  $\Rightarrow$  Different Benefits

Figure 1: Two focal categories of signaling equilibria. In each, one “active” type of agent engages in a signaling behavior while the other “passive” type does not. To sustain a separating equilibrium, it must be that acting generates a net gain for the active type but a net loss for the passive type. In the first category of equilibrium (sub-figure a), both agent types receive the same benefit from engaging in a behavior. The fact that one type acts but the other does not implies that they incur different costs. In the second (sub-figure b), both types incur the same cost. Here, the fact that only one type acts implies that they expect different benefits.

societies (Henrich et al., 2010), however; anthropologists have applied the concept to make sense of puzzling phenomena across the ethnographic record (Bird and Smith, 2005). Biologists, for their part, have used costly signaling to explain cooperative alert calls (Bergstrom and Lachmann, 2001), ritual fighting (Zahavi, 1975), singing (Mithen, 2006), secondary sex characteristics such as bright coloring (Folstad and Karter, 1992), and a variety of other facts about the animal world.

Although, as the above examples make clear, any costly resource can underwrite

the credibility of a signaling equilibrium, which resources happen to be relatively scarce or desirable will, of course, vary across times, places, and cultures.<sup>2</sup> Mental resources are somewhat unique in this regard: in contrast to material goods, whose production depends on technology, capital, and other variable factors, the “one person, one brain” principle—*i.e.*, fact that every individual is endowed a limited set of cognitive capacities (such as working memory and cognitive control; Miller, 1956; Shenhav et al., 2017) that cannot, on a physical level, be increased—ensures that the supply of human attention is essentially fixed (Loewenstein and Wojtowicz, 2023).<sup>3</sup>

Many instances of social signaling rely on the fact that mental resources are scarce and therefore costly to expend. Mental activity entails both direct, metabolic costs (*e.g.*, firing neurons consumes oxygenated blood glucose and generates toxic metabolites; Attwell and Laughlin, 2001) and indirect, opportunity costs (*e.g.*, committing a limited-capacity resource such as cognitive control or working memory to one task precludes its simultaneous use for other, potentially valuable tasks; Shenhav et al., 2017). On the one hand, this means that nearly every act of human intelligence incurs real costs; on the other, it means that all such acts have the potential to credibly signal one’s preferences and abilities to others.

The widespread commercial availability of highly capable artificial intelligence, however, has made it so that anyone can now quickly and cheaply simulate behaviors that previously required the concerted application of these mental resources. A simple example, discussed in detail in Section 5.1, is the act of writing a sincere apology—a mentally effortful activity that can now be automated with a simple prompt. By driving the cost of these mental activities to zero, artificial intelligence has the potential to undermine the very possibility of using signaling as a vehicle for establishing mental proofs across an increasingly wide variety of domains.

Signaling theory predicts that the extent of these effects will depend upon the cost structure supporting the existing communication equilibrium and the manner in which artificial intelligence alters these costs. Consider the case depicted in Figure 1a, where a signaling behavior distinguishes between types who incur high and low cognitive costs. If artificial intelligence enables the (formerly) high-cost types to effect the action at the same cost as the (formerly) low-cost types, then it will no longer differentiate types and its signaling power will collapse. However, if the introduction of artificial intelligence halves both costs, then some difference will remain and the behavior will retain some signaling power. Even in this latter case, the introduction of artificial intelligence will constrain the magnitude of benefits that can be made

---

<sup>3</sup>Bird and Smith (2005) provide examples that range from particularly labor-intensive yams to turtle meat.

<sup>3</sup>If anything, mental resources have become relatively more scarce over time in modern industrial economies as technology drives the cost of producing both material commodities (Nordhaus, 1996) and information (Simon, 1971) toward zero.

contingent upon the signal.

## 4 Mechanism Two: Proofs of Knowledge

Many forms of knowledge are difficult or even impossible to articulate (so-called “tacit” knowledge; Polanyi, 1958; Miton and DeDeo, 2022). In social contexts, this has the important implication that one cannot reveal certain relevant states of knowledge to others through mere recitation. This, in turn, can create barriers to cooperation, as many advantage social arrangements hinge on people possessing specific abilities (*e.g.*, a contractual employment relationship requiring domain-specific expertise).

Fortunately, people can often prove they possess the underlying ability indirectly, by engaging in behaviors that would be prohibitively costly or impossible to perform without it. Beginning with Goldwasser et al. (1985), a sub-branch of cryptography has formally studied the properties of protocols that enable one party (a “prover”) to indirectly convince another party (a “verifier”) that they possess a particular piece of knowledge (*e.g.*, the solution to an intractable problem) or ability (*e.g.*, can compute the values a specific function), without revealing much, if any, additional information beyond the mere possession of that knowledge or ability itself.

Research into such protocols has clarified their essential structural features. According to framework originally laid out by Goldwasser et al. (1985), proof of knowledge protocols must satisfy two conditions: “completeness” (a sincere prover can always convince the verifier) and “validity” (a disingenuous prover cannot convince the verifier, except with arbitrarily small probability).

Complete and valid protocols rely on the fact that knowledge enables us to do new things. Some actions would be improbable for someone who lacked the knowledge in question, and are therefore diagnostic of possessing it. For example, solving certain types of computationally intractable problems immediately entails solutions to a variety of other, structurally-related problems. Demonstrating capacity on these related problems convinces the verifier that the prover must, indeed, possess a solution to the original problem (Goldreich et al., 1991; Blum, 1986). These insights have enabled a variety of applications, such as confidential voting (Groth, 2005), privacy-preserving machine learning (Minelli, 2018), decentralized payments (Sasson et al., 2014), and secure smart cards (Schnorr, 1990).

As Goldwasser et al. (1985) point out, interactivity can afford distinct advantages: “Writing, down a proof that can be checked by everybody without interaction is a much harder task. In some sense, because one has to answer in advance all possible questions” (pg. 292).<sup>4</sup> In the social domain, interactive proofs are, of course, familiar from standardized tests, where a short but probing evaluation of specific problems

enables an examiner to verify a more general capacity. They are also essential to classroom instruction, where students “may ask questions at crucial points of the argument and receive answers” (*Op. cit.*).

Proof of knowledge exchanges extend well beyond the highly structured tests found in licensing exams and technical interviews, however: as we detail below, similar principles enable a sincere apology to certify that a wrongdoer does, indeed, possess a more elaborate mental model of their friend’s needs and goals than the original *faux pas* suggested.

The validity of any such protocol, however, rests upon assumptions about the correlation between one’s observable capabilities and possession of an underlying base of enabling knowledge. It is just this correlation that generative artificial intelligence disrupts: with access to a well-prompted machine, one can simulate many behaviors that would have previously required a far wider spectrum of personal abilities than just crafting a good prompt.

## 5 Examples

### 5.1 A Simple Example: Sincere Apology

In practice, many acts of mental proof in the social domain combine features of both costly signaling and proof of knowledge protocols to certify multiple mental facts simultaneously.

Consider the “sincere apology,” a vital repair mechanism for human social relationships (Bachman and Guerrero, 2006). Although every transgressor would benefit from being forgiven, victims often have good reason not to: inductive evidence that the transgressor may behave poorly in the future. The very fact that apology follows a transgression means that contrary evidence must be communicated at a moment of low trust, hence the technology of mental proof is often vital to establishing one’s credibility.

To see why, consider the four points that Lazare (2005) lists as necessary for an effective apology: (1) identify both offending and offended parties; (2) acknowledge the incident in detail; (3) recognize the harms done; and (4) affirm that the behavior in question violated a social norm. To satisfy these points, the person apologizing must have a good mental model of the events in question and how they affected the injured party: in this sense, it is part of a proof of knowledge exchange. An explicit recitation of events not only establishes that the offending party knows what went wrong, but also generates common knowledge (Chwe, 2013) between both parties

---

<sup>4</sup>Not all protocols, however, require interaction Blum et al. (2019).



about the record of events and their significance relative to the injured party’s needs and goals.

Thinking through a situation in sufficient depth to have explored its various causes and consequences typically requires a significant, concerted expenditure of mental resources on cognitive operations such as mental simulation and perspective taking. One can demonstrate these efforts by making an apology particularly detailed and extensive, pairing it with a “symbolic gesture,” or by presenting novel insight into the situation. No matter the specific method employed, an effortful apology also serves as a costly signal that the offending party values the relationship enough to have invested in the process of rethinking and atoning for their actions.

A successful combination of proof of knowledge and costly signaling can provide strong evidence against the inductive hypothesis of continued transgression by proving that the offending party both: (1) understands the injured party’s needs (proof of knowledge); and (2) values the relationship enough not to repeat the violation (costly signaling). Both parts are required: a heartfelt apology that misunderstands the transgression fails, no matter the effort, while an apology that appears easy and off-the-cuff fails, no matter how accurately it describes the situation.

Recognizing apology as an act of mental proof helps explain why a “ChatGPT apology”<sup>5</sup> written by artificial intelligence does not, to many, count as an apology at all. Recent experiments bear this out: Glikson and Asscher (2023) find that people rate an apology as less sincere and are less likely to forgive when it was written using the heavy use of artificial intelligence tools. Interestingly, Glikson and Asscher also found that the light use of tools (*e.g.*, to correct spelling or grammar errors) incurred no authenticity penalty, presumably because they did not meaningfully reduce perceived cognitive investment in the apology and, therefore, its capacity to deliver mental proof.

## 5.2 A Complex Example: Social Proof

Signaling and proof of knowledge reveal information about the preferences and proficiencies of the agents who engage them. By their very nature, therefore, mental proofs can only directly certify *personal* facts about one’s own mind.

In certain contexts, however, individual facts combine to establish social facts. If a representative sample of people individually prove a mental fact, then an observer can statistically infer that it holds in a broader population. This is especially true in situations where coordination is unlikely or impossible.

Consider, for example, the powerful impact that encountering someone who can speak passionately about a highly specific interest that you, yourself happen to share.

---

<sup>5</sup><https://abcnews.go.com/Business/chatgpt-wedding-vows-eulogies-stokes-dispute-authenticity/story?>

Such an encounter proves not just that you are not alone, but that there are many people like you. Groups establish not just identities, but common knowledge of their size and scope through mental proofs that involve hard-to-acquire knowledge of niche social facts, styles of speech, and shared beliefs.

The capacity of mental proofs to certify both personal and social facts provides one explanation for the emergence of the baroque jargon and world-views in counter-normative communities, such as those devoted to conspiracy theories and extremist ideologies. While some jargon can serve as a shibboleth—a hard-to-forge marker of identity—other jargon is esoteric: learnable, certainly, but only with significant effort by true aficionados (Perry and DeDeo, 2021).

The role of mental proof in establishing social facts provides insight into recent work on artificially-enabled “disinformation” online. Naive observers, for example, who encounter increasing numbers of conspiracy-minded opinions online will naturally take that as evidence for the more general validity of the underlying belief. Such an effect, however, works only as long as the observers remain naive: once it becomes general knowledge that the behavior can be simulated with zero cost, we expect the opposite effect: the discounting of sincerely-professed beliefs and the reach of social movements that advance them.

## 6 When Mental Proof Matters Most

Mental proofs are a pervasive features of social life; Table 1 lists a few examples and describes how Generative AI undermines their social impact. Mental proofs are especially valuable in situations where both (1) deception is potentially profitable and (2) people cannot establish credibility in other ways. They are, therefore, particularly important before trust has been established at the beginning of a romantic, personal, or professional relationship.

Mental proofs also play an important role when when interactions are anonymous or infrequent, as happens in much online communication, or of such high stakes that a counter-party may be willing to sacrifice their reputation for an advantageous outcome. Finally, mental proofs are especially useful for communicating claims that are difficult to verify or articulate (*e.g.*, abstract subjective claims such as caring or understanding).

Mental proofs are powerful, decentralized tools for building cooperation, but they are also, by their nature, costly; for this reason, societies often build more efficient mechanisms around, or on top of, them. If a society makes it possible for trading partners to interact over long periods of time, the accumulation of evidence can make costly proofs less necessary. Trusted reputation systems can ease interactions between strangers by providing a record of past acts of mental proof.

Activity	Mental fact demonstrated	AI Harm to Costly Signaling	AI Harm to Proof of Knowledge
Apology (Section 5.1)	“They have contemplated how their actions affected me and now understand why they shouldn’t behave similarly in the future.”	“They may have used an LLM to write this. If so, they do not actually care enough about me to think through the negative effects their actions had on me.”	“They may not actually understand why the actions harmed me or possess specific background knowledge necessary to avoid harming me in the future.”
Subcultural Membership (Section 5.2)	“Making an obscure reference shows that they have spent a lot of time with the artifacts of my subculture and are highly devoted to its associated values.”	“They may be using an LLM to imitate the communication style of our community. They may have no investment in my subculture and are unlikely to share my values.”	“I cannot be sure if they know other facts or references associated with my subculture.”
Cold Email	“They have demonstrated specific interest in working for my company through extensive internet research and can write intelligently about our industry.”	“They may have used an LLM to write similar messages to a wide array of companies and have no particular interest in my company or the problems it handles.”	“They may have used an LLM to simulate their expertise and do not possess the relevant underlying skills.”
Crisis Hotline Volunteer	“Someone values me enough to have listened to my troubles and responded with encouragement.”	“They may have just fed my message into an LLM, which will return a generic responses regardless of what one says. I am truly alone.”	—
Technical Interview	“By demonstrating obscure knowledge about an aspect of our domain, they show they are familiar enough with it that we will benefit from hiring them.”	—	“They may have gotten this fact from an LLM, without actually possessing any relevant background knowledge. They are not necessarily a good hire.”

Table 1: Examples of how generative AI can undermine the two principle mechanisms of mental proof: costly signaling and proof of knowledge.

Mental proofs are most important in contexts where these supervening mechanisms are not present, for example because institutions are, or have become, weak or misaligned. Students for countries with failed credentialing systems, for example, will be asked to provide more mental proofs compared to others. Potential romantic partners may demand more mental proofs from each other in cultures that can not, or do not, enforce standards of care.

## 7 When AI Damages Mental Proof

Mental proofs are a powerful tool for establishing both individual relationships and the common knowledge necessary for effective group action. Conversely, the weakening of mental proof can significantly stunt people’s ability to form new interpersonal relationships and cooperative initiatives. This suggests that the vibrancy of mental proof is a precursor to the general health of social, familial, political, and economic life in a community.

As the example of sincere apology makes clear, mental proof also play an important role in repairing and maintaining close ties. The general erosion of mental proof therefore also has the effect of diminishing the psychological and social benefits such relationships confer. Feeling that one is understood and cared about by another is a deep psychological priority (Cahn, 1990; Oishi et al., 2010; Reis et al., 2000; Lun et al., 2008), separate from the material benefits that care might bring.

Our discussion of social proof also showed how mental proofs help people establish collective interests, beliefs, and capacities. These proofs provide more than just the psychological benefits of knowing one is not alone. Collective action matters—the creation and maintenance of groups is a basic feature of civil society, in general, and the success of democratic government, in particular (Putnam, 1994).

## 8 Implications

Artificial intelligence’s deleterious effects on both costly signaling and proofs of knowledge can be prevented if people can clearly delineate between communicative acts undertaken with and without the assistance of such tools. In the United States, this is a stated goal of a 2023 Executive Order (No. 14110, Section 4.5),<sup>6</sup> and Jain et al. (2023) detail a variety of strategies to establish and maintain such distinctions, which they label the *contextual confidence* of communication. Clarifying the difference between AI and human content would enable us to reap the benefits of automation (where advantageous) while preserving the capacity of humans to harness the benefits of mental proof (where important).

---

<sup>6</sup><https://www.whitehouse.gov/briefing-room/presidential-actions/2023/10/30/executive-order-on-the->

Our analysis also suggests that the economic benefits, in terms of reduced labor costs, of automating human mental effort can backfire. This may be especially true for jobs where care and understanding are paramount. Consider a patient in talk therapy for a difficult-to-treat condition. Generative artificial intelligence may help the therapist by spotting patterns and dynamically recommending better strategies of engagement, but we do not yet understand how this might undermine the therapeutic alliance (Tal et al., 2023; Zetzel, 1956). Clarifying the work being done by mental proof may help guide efforts to surgically target aspects of these jobs that can benefit from automation without damaging their core efficacy.

Our work has focused on the beneficial effects of mental proof. As Spence (1973) pointed out when introducing the concept, however, the existence of signaling equilibria can be socially costly and may create disadvantages for some participants. This suggests, in turn, that in some cases—which ones remains a subject for further research—outcomes may be improved when Generative AI destroys an equilibrium previously supported by mental proof. The conditions under which this is a net gain to the participants, either individually, or collectively, depends sensitively on the mechanisms that emerge to take its place.

Mental proof is most prominent in low-trust environments; consequently, its degradation stands to disproportionately impact those who already lack strong institutional support. Consider, for example, a student in the developing world without access to a functional accreditation system. In the past, a well-crafted, thoughtful e-mail might well serve to open doors to informal networks of mentorship and training: such a gesture provided both proof of abilities and of the necessary internal motivation that would lead a busy, but sympathetic, professor to take note. Such avenues to advancement are closed, however, once equivalent texts can be produced with a minimal prompt. A student without social capital is hurt by this vitiation of mental proof, while another, who comes with institutional backing, has other avenues to establish their credibility.

This suggests that expanding access to institutions that establish and sustain trust will be even more valuable for maintaining open and equitable societies as artificial intelligence advances. Novel protocols that bootstrap existing sources of trust (*e.g.*, Weyl et al., 2022) could help counteract the inevitable changes that artificial intelligence will bring to the cost of thought and, through it, the structure of mental proof.

## References

Ella Glikson and Omri Asscher. AI-mediated apology in a multilingual work context: Implications for perceived authenticity and willingness to forgive. *Computers in*

- Human Behavior*, 140:107592, 2023.
- Tira Nur Fitria. Artificial intelligence (AI) technology in OpenAI ChatGPT application: A review of ChatGPT in writing English essay. In *ELT Forum: Journal of English Language Teaching*, volume 12, pages 44–58, 2023.
- Peter Cardon, Carolin Fleischmann, Jolanta Aritz, Minna Logemann, and Jeanette Heidewald. The challenges and opportunities of ai-assisted writing: Developing ai literacy for the ai age. *Business and Professional Communication Quarterly*, 86(3): 257–295, 2023.
- Yihan Wu and Ryan M Kelly. Online dating meets artificial intelligence: How the perception of algorithmically generated profile text impacts attractiveness and trust. In *Proceedings of the 32nd Australian Conference on Human-Computer Interaction*, pages 444–453, 2020.
- Tammy LaGorce. Need to write your vows? A.I. can help. 2023. URL <https://www.nytimes.com/2023/03/03/fashion/weddings/chatbot-wedding-vows-chatgpt-ai>
- Ella Glikson and Anita Williams Woolley. Human trust in artificial intelligence: Review of empirical research. *Academy of Management Annals*, 14(2):627–660, 2020.
- Shrey Jain, Zoë Hitzig, and Pamela Mishkin. Contextual confidence and generative ai. *arXiv preprint arXiv:2311.01193*, 2023.
- David Leslie. Understanding artificial intelligence ethics and safety. *arXiv preprint arXiv:1906.05684*, 2019.
- Danielle Allen and E Glen Weyl. The real dangers of generative ai. *Journal of Democracy*, 35(1):147–162, 2024.
- Andreas Jungherr. Artificial intelligence and democracy: A conceptual framework. *Social media+ society*, 9(3):20563051231186353, 2023.
- Alec Tyson. Growing public concern about the role of artificial intelligence in daily life, Aug 2023. URL <https://www.pewresearch.org/short-reads/2023/08/28/growing-public-concern-about-the>
- Irene Solaiman, Zeerak Talat, William Agnew, Lama Ahmad, Dylan Baker, Su Lin Blodgett, Hal Daumé III, Jesse Dodge, Ellie Evans, Sara Hooker, et al. Evaluating the social impact of generative ai systems in systems and society. *arXiv preprint arXiv:2306.05949*, 2023.

- Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, et al. Ethical and social risks of harm from language models. *arXiv preprint arXiv:2112.04359*, 2021.
- Yisroel Mirsky and Wenke Lee. The creation and detection of deepfakes: A survey. *ACM Computing Surveys (CSUR)*, 54(1):1–41, 2021.
- Michael Tomasello, Malinda Carpenter, Josep Call, Tanya Behne, and Henrike Moll. Understanding and sharing intentions: The origins of cultural cognition. *Behavioral and brain sciences*, 28(5):675–691, 2005.
- Margaret Gilbert. Walking together: A paradigmatic social phenomenon. *MidWest studies in philosophy*, 15(1):1–14, 1990.
- Michael E Bratman. Shared cooperative activity. *The philosophical review*, 101(2):327–341, 1992.
- Joseph Farrell. Cheap talk, coordination, and entry. *The RAND Journal of Economics*, pages 34–39, 1987.
- Ernst Fehr. On the economics and biology of trust. *Journal of the european economic association*, 7(2-3):235–266, 2009.
- Jon Elster. Social norms and economic theory. *Journal of economic perspectives*, 3(4):99–117, 1989.
- Douglass C North. Institutions. *Journal of economic perspectives*, 5(1):97–112, 1991.
- Michael Spence. Job market signaling. *The Quarterly Journal of Economics*, 87(3):355–374, 1973.
- Amotz Zahavi. Mate selection—a selection for a handicap. *Journal of theoretical Biology*, 53(1):205–214, 1975.
- John Maynard Smith and David Harper. *Animal signals*. Oxford University Press, 2003.
- Claire D Fitzgibbon and John H Fanshawe. Stotting in thomson’s gazelles: an honest signal of condition. *Behavioral Ecology and Sociobiology*, 23:69–74, 1988.
- Tim M Caro. Ungulate antipredator behaviour: preliminary and comparative data from african bovids. *Behaviour*, 128(3-4):189–228, 1994.

- Paul Milgrom and John Roberts. Price and advertising signals of product quality. *Journal of political economy*, 94(4):796–821, 1986.
- Beiting Cheng, Ioannis Ioannou, and George Serafeim. Corporate social responsibility and access to finance. *Strategic management journal*, 35(1):1–23, 2014.
- Caroline Flammer. Corporate green bonds. *Journal of financial economics*, 142(2):499–516, 2021.
- Gerrit KC Ahlers, Douglas Cumming, Christina Günther, and Denis Schweizer. Signaling in equity crowdfunding. *Entrepreneurship theory and practice*, 39(4):955–980, 2015.
- Joseph Henrich, Steven J Heine, and Ara Norenzayan. Most people are not WEIRD. *Nature*, 466(7302):29–29, 2010.
- Rebecca Bird and Eric Smith. Signaling theory, strategic interaction, and symbolic capital. *Current anthropology*, 46(2):221–248, 2005.
- Carl T Bergstrom and Michael Lachmann. Alarm calls as costly signals of antipredator vigilance: the watchful babbler game. *Animal behaviour*, 61(3):535–543, 2001.
- Steven J Mithen. *The singing Neanderthals: The origins of music, language, mind, and body*. Harvard University Press, 2006.
- Ivar Folstad and Andrew John Karter. Parasites, bright males, and the immunocompetence handicap. *The American Naturalist*, 139(3):603–622, 1992.
- George A Miller. The magical number seven, plus or minus two: Some limits on our capacity for processing information. *Psychological review*, 63(2):81, 1956.
- Amitai Shenhav, Sebastian Musslick, Falk Lieder, Wouter Kool, Thomas L Griffiths, Jonathan D Cohen, and Matthew M Botvinick. Toward a rational and mechanistic account of mental effort. *Annual review of neuroscience*, 40:99–124, 2017.
- George Loewenstein and Zachary Wojtowicz. The economics of attention. *Available at SSRN 4368304*, 2023.
- William D Nordhaus. Do real-output and real-wage measures capture reality? the history of lighting suggests not. In *The economics of new goods*, pages 27–70. University of Chicago Press, 1996.
- Herbert A. Simon. Designing organizations for an information-rich world. 1971.



- David Attwell and Simon B Laughlin. An energy budget for signaling in the grey matter of the brain. *Journal of Cerebral Blood Flow & Metabolism*, 21(10):1133–1145, 2001.
- Michael Polanyi. Personal knowledge: Towards a post-critical philosophy. 1958.
- Helena Miton and Simon DeDeo. The cultural transmission of tacit knowledge. *Journal of the Royal Society Interface*, 19(195):20220238, 2022.
- S Goldwasser, S Micali, and C Rackoff. The knowledge complexity of interactive proof-systems. In *Proceedings of the seventeenth annual ACM symposium on Theory of computing*, pages 291–304, 1985.
- Oded Goldreich, Silvio Micali, and Avi Wigderson. Proofs that yield nothing but their validity or all languages in np have zero-knowledge proof systems. *Journal of the ACM (JACM)*, 38(3):690–728, 1991.
- Manuel Blum. How to prove a theorem so no one else can claim it. In *Proceedings of the International Congress of Mathematicians*, volume 1, page 2. Citeseer, 1986.
- Jens Groth. Non-interactive zero-knowledge arguments for voting. In *Applied Cryptography and Network Security: Third International Conference, ACNS 2005, New York, NY, USA, June 7-10, 2005. Proceedings 3*, pages 467–482. Springer, 2005.
- Michele Minelli. *Fully homomorphic encryption for machine learning*. PhD thesis, Université Paris sciences et lettres, 2018.
- Eli Ben Sasson, Alessandro Chiesa, Christina Garman, Matthew Green, Ian Miers, Eran Tromer, and Madars Virza. Zerocash: Decentralized anonymous payments from bitcoin. In *2014 IEEE symposium on security and privacy*, pages 459–474. IEEE, 2014.
- Claus-Peter Schnorr. Efficient identification and signatures for smart cards. In *Advances in Cryptology—CRYPTO’89 Proceedings 9*, pages 239–252. Springer, 1990.
- Manuel Blum, Paul Feldman, and Silvio Micali. Non-interactive zero-knowledge and its applications. In *Providing Sound Foundations for Cryptography: On the Work of Shafi Goldwasser and Silvio Micali*, pages 329–349. 2019.
- Guy Foster Bachman and Laura K Guerrero. Forgiveness, apology, and communicative responses to hurtful events. *Communication reports*, 19(1):45–56, 2006.
- Aaron Lazare. *On apology*. Oxford University Press, 2005.

- Michael Suk-Young Chwe. *Rational ritual: Culture, coordination, and common knowledge*. Princeton University Press, 2013.
- Chloe Perry and Simon DeDeo. The cognitive science of extremist ideologies online. *arXiv*, 2110.00626, 2021.
- Dudley D Cahn. Perceived understanding and interpersonal relationships. *Journal of Social and Personal Relationships*, 7(2):231–244, 1990.
- Shigehiro Oishi, Margarita Krochik, and Sharon Akimoto. Felt understanding as a bridge between close relationships and subjective well-being: Antecedents and consequences across individuals and cultures. *Social and Personality Psychology Compass*, 4(6):403–416, 2010.
- Harry T Reis, Kennon M Sheldon, Shelly L Gable, Joseph Roscoe, and Richard M Ryan. Daily well-being: The role of autonomy, competence, and relatedness. *Personality and social psychology bulletin*, 26(4):419–435, 2000.
- Janetta Lun, Selin Kesebir, and Shigehiro Oishi. On feeling understood and feeling well: The role of interdependence. *Journal of Research in Personality*, 42(6):1623–1628, 2008.
- Robert D Putnam. *Making democracy work: Civic traditions in modern Italy*. Princeton University Press, 1994.
- Amir Tal, Zohar Elyoseph, Yuval Haber, Tal Angert, Tamar Gur, Tomer Simon, and Oren Asman. The artificial third: utilizing chatgpt in mental health. *The American Journal of Bioethics*, 23(10):74–77, 2023.
- Elizabeth R. Zetzel. An approach to the relation between concept and content in psychoanalytic theory. *The Psychoanalytic Study of the Child*, 11(1):99–121, 1956. doi: 10.1080/00797308.1956.11822784.
- E Glen Weyl, Puja Ohlhaver, and Vitalik Buterin. Decentralized society: Finding web3’s soul. *Available at SSRN 4105763*, 2022.