Personhood credentials: Artificial intelligence and the value of privacy-preserving tools to distinguish who is real online

Steven Adler,*^{†1} Zoë Hitzig,*^{†1,2} Shrey Jain,*^{†3} Catherine Brewer,*⁴ Wayne Chang,*⁵ Renée DiResta,*²⁵ Eddy Lazzarin,*⁶ Sean McGregor,*⁷ Wendy Seltzer,*⁸ Divya Siddarth,*⁹ Nouran Soliman,*¹⁰ Tobin South,*¹⁰ Connor Spelliscy,*¹¹ Manu Sporny,*¹² Varya Srivastava,*⁴ John Bailey,¹³ Brian Christian,⁴ Andrew Critch,¹⁴ Ronnie Falcon,¹⁵ Heather Flanagan,²⁵ Kim Hamilton Duffy,¹⁶ Eric Ho,¹⁷ Claire R. Leibowicz,¹⁸ Srikanth Nadhamuni,¹⁹ Alan Z. Rozenshtein,²⁰ David Schnurr,¹ Evan Shapiro,²¹ Lacey Strahm,¹⁵ Andrew Trask,^{4,15} Zoe Weinberg,²² Cedric Whitney,²³ Tom Zick²⁴

¹OpenAI, ²Harvard Society of Fellows, ³Microsoft, ⁴University of Oxford, ⁵SpruceID, ⁶a16z crypto,
 ⁷UL Research Institutes, ⁸Tucows, ⁹Collective Intelligence Project, ¹⁰Massachusetts Institute of Technology,
 ¹¹Decentralization Research Center, ¹²Digital Bazaar, ¹³American Enterprise Institute,
 ¹⁴Center for Human-Compatible AI, University of California, Berkeley, ¹⁵OpenMined,
 ¹⁶Decentralized Identity Foundation, ¹⁷Goodfire, ¹⁸Partnership on AI, ¹⁹eGovernments Foundation,
 ²⁰University of Minnesota Law School, ²¹Mina Foundation, ²²ex/ante, ²³School of Information, University of California, Berkeley,
 ²⁴Berkman Klein Center for Internet & Society, Harvard University, ²⁵Independent Researcher

August 2024

Abstract

Anonymity is an important principle online. However, malicious actors have long used misleading identities to conduct fraud, spread disinformation, and carry out other deceptive schemes. With the advent of increasingly capable AI, bad actors can amplify the potential scale and effectiveness of their operations, intensifying the challenge of balancing anonymity and trustworthiness online. In this paper, we analyze the value of a new tool to address this challenge: "personhood credentials" (PHCs), digital credentials that empower users to demonstrate that they are real people—not AIs—to online services, without disclosing any personal information. Such credentials can be issued by a range of trusted institutions—governments or otherwise. A PHC system, according to our definition, could be local or global, and does not need to be biometrics-based. Two trends in AI contribute to the urgency of the challenge: Al's increasing indistinguishability from people online (i.e., lifelike content and avatars, agentic activity), and AI's increasing scalability (i.e., cost-effectiveness, accessibility). Drawing on a long history of research into anonymous credentials and "proof-of-personhood" systems, personhood credentials give people a way to signal their trustworthiness on online platforms, and offer service providers new tools for reducing misuse by bad actors. In contrast, existing countermeasures to automated deception—such as CAPTCHAs—are inadequate against sophisticated AI, while stringent identity verification solutions are insufficiently private for many use-cases. After surveying the benefits of personhood credentials, we also examine deployment risks and design challenges. We conclude with actionable next steps for policymakers, technologists, and standards bodies to consider in consultation with the public.

[†] Indicates the corresponding authors: Steven Adler (steven@openai.com), Zoë Hitzig (zhitzig@openai.com), and Shrey Jain (shrey-jain@microsoft.com).

^{*} Denotes primary authors, who contributed most significantly to the direction and content of the paper. Besides corresponding authors, all other authors are listed in alphabetical order.

Executive Summary

Malicious actors have long used misleading identities to deceive others online. They carry out fraud, cyberattacks, and disinformation campaigns from multiple online aliases, email addresses, and phone numbers. Historically, such deception has sometimes seemed an unfortunate but necessary cost of preserving the Internet's commitments to privacy and unrestricted access. But highly capable AI systems may change the landscape: There is a substantial risk that, without further mitigations, deceptive AI-powered activity could overwhelm the Internet. To uphold user privacy while protecting against AI-powered deception, new countermeasures are needed.

With access to increasingly capable AI, malicious actors can potentially orchestrate more effective deceptive schemes. Two trends contribute to these schemes' potential impact:

- 1. **Indistinguishability.** Distinguishing AI-powered users on the Internet is becoming increasingly difficult, as AI advances in its ability to:
 - <u>Generate human-like content</u> that expresses human-like experiences or points of view (e.g., "Here is what I thought of that speech").
 - <u>Create human-like avatars</u> through photos, videos, and audio (e.g., simulating a real-looking person on a video chat).
 - <u>Take human-like actions</u> across the Internet (e.g., browsing websites like an ordinary user, making sophisticated plans to achieve goals they are given, solving CAPTCHAs when challenged).
- 2. Scalability. AI-powered deception by malicious actors is increasingly scalable because of:
 - Decreasing costs at all capability levels.
 - Increasing accessibility, for example, via open-weights deployments through which scaled misuse is less preventable.

Taken together, these two trends suggest that AI may help to make deceptive activity more convincing (through increased indistinguishability) and easier to carry out (through increased scalability).

We identify one promising solution to pervasive deception on the Internet, building off decades of research in cryptography and experimentation in online communities: personhood credentials (hereafter referred to as PHCs). Such a credential empowers its holder to demonstrate to providers of digital services that they are a person without revealing anything further. Building on related concepts like proof-of-personhood and anonymous credentials, these credentials can be stored digitally on holders' devices and verified through zero-knowledge proofs. Importantly, such proofs do not reveal the individual's specific credential (nor any aspects of their identity).

To counter scalable deception while maintaining user privacy, PHC systems must meet two foundational requirements:

- 1. Credential limits: The issuer of a PHC gives at most one credential to an eligible person.
- 2. Unlinkable pseudonymity: PHCs let a user interact with services anonymously through a service-specific pseudonym; the user's digital activity is untraceable by the issuer and unlinkable across service providers, even if service providers and issuers collude.

These two properties equip service providers with the option to offer services on a per-person basis, and to prevent the return of users who violate the service's rules. An anonymous forum, for instance, could offer a single verified account to each credential holder. Unlinkable pseudonymity helps them achieve this because it prevents one person from using the same PHC to sign up twice, even without ever identifying the user. The issuer's credential limit gives them high confidence that the same user cannot easily circumvent the limit by using many PHCs to make many different accounts.

There are many effective ways to design a PHC system, and various organizations—governmental or otherwise—can serve as issuers. In one possible implementation, states could offer a PHC to any holder of their state's tax identification number; a PHC system, according to our definition, could be local or global, and

does not need to be based in biometrics. Having multiple trusted PHC issuers within a single ecosystem promotes choice—people can select into systems built on their preferred root of trust (government IDs, social graphs, biometrics) and that offer affordances that best align with their preferences. This approach reduces the risks associated with a single centralized issuer while still preserving the ecosystem's integrity by limiting the total number of credentials. Note that this paper does not advocate for or against any specific PHC system design; instead, it aims to establish the value of PHCs in general while highlighting challenges that must be taken into account in any design.

PHCs are not forgeable by AI systems, and it is difficult for malicious actors to obtain many of them. By combining verification techniques that have an offline component (e.g., appearing in-person, validating a physical document) and secure cryptography, these credentials are issued only to people and cannot be convincingly faked thereafter. They therefore help to counter the problem of indistinguishability by creating a credential only people can acquire, and help to counter the problem of scalability by enabling per-credential rate limits on activities.

PHCs give digital services a tool to reduce the efficacy and prevalence of deception, especially in the form of:

- 1. Sockpuppets: deceptive actors purporting to be "people" that do not actually exist.
- 2. <u>Bot attacks</u>: networks of bots controlled by malicious actors to carry out automated abuse (e.g., breaking site rules and evading suspension by creating new accounts).
- 3. Misleading agents: AI agents misrepresenting whose goals they serve.

PHCs offer people a tool to credibly signal that they are a real person operating an authentic account, without conveying their identity. PHCs also help service providers spot deceptive accounts, which may lack such a signal.

PHCs improve on and complement existing approaches to countering AI-powered deception online. For example, the following approaches are often not robust to highly capable AI, not inclusive, and/or not privacy-preserving:

- 1. Behavioral filters, e.g., CAPTCHAs, JavaScript browser challenges, anomaly detection.
- 2. Economic barriers, e.g., paid subscriptions, credit card verification.
- 3. AI content detection, e.g., watermarking, fingerprinting, metadata provenance.
- 4. Appearance- and document-based verification, e.g., selfie checks with ID, live video calls.
- 5. Digital and hardware identifiers, e.g., phone numbers, email addresses, hardware security keys.

To achieve their benefits, **PHC systems must be designed and implemented with care.** We discuss four areas in which PHCs' impacts must be carefully managed:

- 1. Equitable access to digital services that use PHCs.
- 2. Free expression supported by confidence in the privacy of PHCs.
- 3. Checks on power of service providers and PHC issuers.
- 4. Robustness to attack and error by different actors in the PHC ecosystem.

In close collaboration with the public, we encourage governments, technologists, and standards bodies to invest in the development, piloting, and adoption of personhood credentials as a key tool in addressing scalable deception online:

- 1. <u>Invest in development</u> and piloting of personhood credentialing systems. <u>e.g.</u>, <u>explore building PHCs incrementally atop existing credentials such as digital driver's licenses.</u>
- 2. Encourage adoption of personhood credentials.

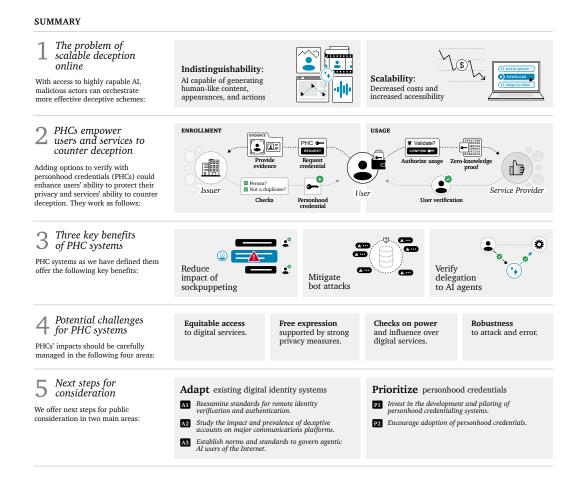
 e.g., determine services for which PHCs ought to be substitutable for ID verification.

It is also important that these groups accelerate their preparations for AI's impact more generally by adapting existing digital systems:

- 1. Reexamine standards for remote identity verification and authentication.
 e.g., reconsider confidence in selfie-based identity verification, absent supplemental factors to reduce AI-enabled spoofing.
- 2. <u>Study</u> the impact and prevalence of deceptive accounts across major communications platforms. e.g., develop standardized methods for measuring the prevalence of fake accounts on social media.
- 3. <u>Establish norms and standards</u> to govern agentic AI users of the Internet. e.g., explore new forms of trust infrastructure for AI agents, akin to HTTPS for websites.

Readers primarily interested in ideas for next steps may refer directly to Section 5 for further detail.

We are concerned that the Internet is inadequately prepared for the challenges highly capable AI may pose. Without proactive initiatives involving the public, governments, technologists, and standards bodies, there is a significant risk that digital institutions will be unprepared for a time when AI-powered agents, including those leveraged by malicious actors, overwhelm other activity online. Lacking better alternatives, institutions might resort to privacy-violating methods for rooting out scaled deception, like creating digital identification systems that (intentionally or unintentionally) link a person's legal identity with a complete record of their digital activity. By contrast, personhood credentials have the potential to reduce deceptive activity while preserving privacy—giving people and services the tools to signal and sustain trustworthiness online.



Contents

Ex	ecutive Summary	2	
1 1.1 1.2 1.3	2 Trends in AI threaten to make online deception more effective	8	
2 2.1 2.2			
3 3.1 3.2 3.3	2 Mitigate bot attacks	19	
4 4.2 4.3 4.4	Free expression	23 24	
5 5.2 5.2	2 Prioritize personhood credentials		
Ac	Acknowledgments		
Re	References		
Ap	pendices	50	
A	What do we mean by "trustworthy" digital interaction?	50	
В	Relating personhood credentials to CAPTCHAs and synthetic content transparency tools	51	
С	Implementation choices for personhood credentials	55	
D	Ecosystem design and management	60	

1 Introduction

1.1 The Internet has long struggled with deceptive activity

Malicious actors have long used misleading identities to carry out abuse online. For example, attackers manipulate perceptions of public opinion by spreading disinformation through deceptive "sockpuppet" accounts—appearing to represent distinct people and thus lending their claims more credibility.¹ Other deceptive attacks employ automated botnets, executing distributed attacks by posing as a large number of distinct users.² In a wide range of digital systems, the ability to create deceptive and seemingly independent profiles can be exploited by malicious actors.³

Defenders of online freedom of expression sometimes view deception as an unfortunate but necessary cost of preserving users' privacy [32]. Indeed, the Internet's pioneers saw anonymity as a fundamental pillar of privacy and freedom of expression [148]. They built its architecture to let people participate in digital spaces without disclosing their real identities. The benefits from this approach are numerous, and worthy of steadfast protection—for example, anonymity allows people in oppressive regimes to express their opinions without fear of retribution. However, these benefits have come at some cost—one such cost is a lack of accountability for deceptive misuse.⁴

Although deceptive activity has significantly taxed the Internet,⁵ the network remains largely usable for ordinary users. Defenses such as spam filters, IP blocklists, firewalls, and vigilant security analysts have helped to detect and mitigate deceptive attacks.⁶ Furthermore, the attacks themselves have been resource-constrained—their reliance on human labor has kept their scale and quality in check.⁷

The wide availability of increasingly capable AI⁸ may upset this balance. Although bad actors have perpetrated deceptive attacks for decades, actors' increased access to sophisticated and inexpensive AI tools may make their attacks far more effective—harder to distinguish and also more prevalent [112].

The resulting escalation in deceptive activity across social media, public comment systems, and other essential digital services could create substantial challenges for institutions that rely on the Internet.

¹ For an overview of the prevalence of these attacks, see [31]. For an influential early attempt to detect astroturfing on Twitter, see [212].

² A "botnet" originally referred to a network of compromised hardware devices, which could be controlled together. With the advent of cloud computing, botnets can also be run from rented computers in the cloud rather than relying on infected host computers. Botnets are used to carry out DDoS and credential stuffing attacks, to spread spyware, and to propagate scams [65, 237].

³ This strategy of wielding multiple distinct identities to manipulate a system, especially when carried out on peer-to-peer networks, is known in computer security as a "Sybil attack." For the foundational treatment, see [78].

⁴ For an account of how these norms evolved on Usenet, see [278]. A common approach to undesired content sharing on Usenet was to apply social pressure [21].

⁵ It is difficult to precisely measure the tax imposed by malicious activity. One 2012 estimate of the cost to US firms and consumers of online spam alone was \$20 billion [211].

⁶ After being caught executing a particular type of attack, a malicious actor can try again. However, this usually involves additional work, like compromising more accounts or configuring a different proxy to disguise one's IP address. If the actor does not change their tactics, they may be caught again by the same detection methods [105].

⁷ For an overview of the types of labor involved in carrying out influence operations—and how AI may transform these—see [112].

⁸ In this paper, when we refer to "AI," we mean the *combination* of available AI models that can be used together for different purposes. Generally, improvements in AI capabilities come from introducing new models that outperform existing ones in some way or offer better performance relative to their cost.

Under pressure, these institutions may resort to invasive measures for verifying users' identities online, overturning the Internet's longstanding commitment to privacy and civil liberties.⁹

In this paper, we focus on one potential solution to counter AI-powered deception: personhood credentials (hereafter referred to as PHCs), which certify that their holder is a person¹⁰ without revealing anything more about their identity. PHCs help to distinguish people from even the most advanced AI systems by relying on two important deficits of AI. Specifically, AI systems cannot convincingly mimic people offline, and they cannot bypass state-of-the-art cryptographic systems.¹¹ PHCs allow people to interact with different digital service providers in a way that maintains their privacy through unlinkable pseudonymity—the credential holder can have a persistent, privacy-preserving pseudonym with different service providers, and different service providers have no way to link their activity. The issuer of the credential cannot trace their activity either.¹²

OFFLINE COMPONENT Pass as a person in the real-world CRYPTOGRAPHY Forge advanced cryptography CRYPTOGRAPHY Forge advanced cryptography

Things that AI cannot do

Figure 1: Personhood credentials rely on two deficits of AI.

There are many possible issuers of such credentials, with many different methods an issuer can use to ensure that these credentials are available only to people and are not misappropriated by bad actors. For example, states could offer a PHC to each holder of a tax identification number, but thereafter have no way of tracing the uses of the credential.

That personhood credentials ought to exist is not a wholly novel idea. Related credentials are the subject of ongoing initiatives, ¹³ as well as decades of research in security and cryptography. In particular, a

⁹ China's Ministry of Public Security and its Cyberspace Administration recently proposed a national system to associate online activity with individuals' legal identities [251]. The proposal has faced criticism from privacy advocates.

¹⁰ For the purposes of this paper, we use "person" to refer to a human being, though we recognize that some domains use the term "person" more broadly. For example, in law, "a person is any being whom the law regards as capable of rights or duties" [223], which can include "artificial persons" (e.g., corporations) in addition to "natural persons" (i.e., human beings) [107]. In moral philosophy, "personhood" need not apply only to human beings and could include other groups like non-human animals; "personhood" is sometimes associated with having full moral status, though philosophers disagree considerably about which traits are necessary or sufficient for full moral status [59]. Though we use "person" and "personhood" to refer to human beings, this should not be read as implying any particular view on whether future AI systems should be considered moral persons or granted legal personhood (whether for moral or pragmatic reasons).

¹¹ Cryptography relies on computationally hard mathematical problems, such as the factoring of very large numbers. There are not any known methods of efficiently solving certain such problems, whether by a human or an AI system.

¹²While some national digital identity systems also draw on cryptography, they have very different goals and privacy affordances compared to PHCs. For example, Estonia's national ID-card allows citizens to conduct a wide range of activities online (e.g., cryptographically signing important contracts). But these systems are not intended to conceal *which* holder of a card is performing the action [213].

¹³ Our advocacy of personhood credentials aligns with many ongoing initiatives that aim for minimal disclosure identity and authentication systems, e.g., World Wide Web Consortium's (W3C's) Verifiable Credentials [240] and decentralized identifiers [239], European Union Digital Identity's (EUDI's) privacy-preserving digital wallets [88], and a range of standards and implementations of anonymous and attribute-based credentials [1, 126, 90]. One particularly relevant implementation is British Columbia's Person Credential [34, 165].

personhood credential is a type of anonymous credential—first proposed in the 1980s¹⁴—and can be understood as similar to a "proof-of-personhood" system.¹⁵ What is unprecedented is the scale and urgency of the problems these credentials could address, and the harms to trustworthy interactions¹⁶ on the Internet that may occur if unmitigated.

1.2 Trends in AI threaten to make online deception more effective

To better understand the urgent need for personhood credentials, we discuss two notable trends in AI development that may contribute to more effective online deception.¹⁷ These trends are summarized in Figure 2.

First, users powered by AI systems are **increasingly indistinguishable** from people online. For example, AI-powered accounts can be populated with realistic human-like avatars that share high-quality content and that take increasingly autonomous actions.

Second, AI is becoming **increasingly scalable**—both more affordable and accessible, which can help for achieving many benefits but can also enable a larger amount of deception online.

AI trends that could contribute to a rise in misleading activity online

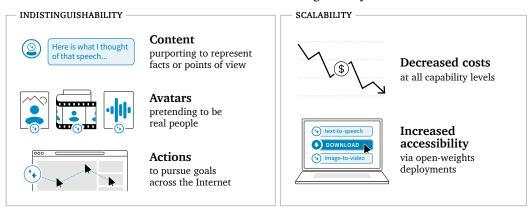


Figure 2: Indistinguishability and scalability could drive an increase in AI-powered deception.

¹⁴ An anonymous credential is a credential that allows its holder to prove some specific statement about themselves without revealing anything further [45, 47, 158]. In the case of a personhood credential, the claim is that its holder is a person. In other cases, one could have an anonymous credential that proves their age without revealing anything further. Personhood credentials largely coincide with "accountable pseudonyms," which aim to link a pseudonymous digital certificate that can be used anonymously online to an offline entity [99]. For a recent systematization of knowledge paper about anonymous credentials, see [118].

¹⁵ The term "proof of personhood" emerged in blockchain communities—used in contrast to the concepts of "proof of work" and "proof of stake"—as a strategy for dealing with Sybil attacks (where one entity subverts a digital system by using many pseudonymous identities to exert undue influence) in permissionless voting systems [29]. Given that "proof of personhood systems" often have specific aims (for instance, to allocate a cryptocurrency or voting rights in a decentralized autonomous organization) and different resulting design goals (e.g., they often aim for global uniqueness, sometimes via biometrics), we use the more generic term "personhood credentials" throughout the paper.

¹⁶ By "trustworthy digital interactions," we mean interactions in which the parties' reasonable expectations are achieved: for instance, a trustworthy online ecosystem is one in which service providers are confident that their services are being used as intended, and users can access digital services free from fears of abuse, attack, and other harms. We further unpack the foundations of this concept in Appendix A. drawing largely on [133].

¹⁷When we refer to "highly capable AI" systems, we mean AI systems that exhibit these properties.

AI systems are increasingly indistinguishable from people online

Malicious actors can employ AI at many stages of their deceptive schemes, drawing on AI's ability to seem like a real person online. This increasing indistinguishability comes from three categories of capabilities:¹⁸

- Generating human-like content that expresses human-like experiences or points of view (e.g., posting commentary like "Here is what I thought of that speech").
- Creating human-like avatars through photos, videos, and audio (e.g., simulating a real-looking person on a video chat).
- Taking human-like actions across the Internet (e.g., browsing websites like an ordinary user, making sophisticated plans to achieve goals they are given, solving CAPTCHAs when challenged).

AI systems are improving in each of these categories. For instance, AI-generated text and images are increasingly difficult to distinguish from human-created content in domains like politics, ¹⁹ art [220], and literature [58]. Already, many people struggle to discern whether they are conversing by text with a human or an AI system in some settings. ²⁰ A growing body of literature explores how AI-generated content might contribute to the efficacy of deceptive schemes, such as through AI's potential abilities to make persuasive arguments [201, 17, 14, 85], tailor persuasion to a specific recipient [162], and roleplay as certain personas [265].

Likewise, AI is improving at creating human-like avatars. AI can enable malicious actors to create realistic high-quality photos and videos of people who do not exist [147, 33, 259, 35]; animate photos of people into videos [233]; take on the appearance of a different person over video chat [276, 123]; and speak with a realistic human voice, whether a generic voice or that of a specific person. These capabilities can be used to spoof identification checks, such as those requiring a selfie with a matching driver's license. In the near future, during a video chat, a user may not be able to discern whether their conversational partner is actually the person they are seeing, a person disguising themselves using AI, or even a complete AI simulation of a real or fictitious person. Simulation of a real or fictitious person.

AI also continues to improve at taking human-like actions across the Internet. For instance, AI systems can solve CAPTCHA puzzles to gain access to Web services that are meant only for humans.²⁴ More generally, AI systems are becoming increasingly agentic: capable of dynamically and independently carrying out

¹⁸ Malicious actors can draw on these capabilities at many stages of their deceptive schemes. For example, a malicious actor could create a pipeline of sockpuppet accounts on a social media site, in which it pairs a realistic-looking AI-generated profile picture (appearance) with a stream of posts that express a particular point-of-view (expressions) and a sequence of how to interact with organic users on the site (actions).

¹⁹ For a study of legislators' response rate to emails drafted by AI compared to those drafted by people—finding only a small difference—see [150].

²⁰ Research suggests that text written by AI systems are hard to distinguish from text written by people, such as in Turing Test conversations [135, 5]. These findings are not wholly new; at 2020's GPT-3 launch, OpenAI reported that participants in a research study had trouble distinguishing between news articles by people versus those generated by AI [36].

²¹ For an overview of voice generation abilities and resultant challenges, see [196]. There are many different types of risks that can emerge with voice cloning [125]. Voice cloning scams, in particular, are rampant and on the rise [4].

²² Some platform verification processes request proof of a government ID card—but even this process cannot restrict entry of AIs working on behalf of bad actors, as AI tools can be used to spoof a realistic-looking selfie of any person holding their supposed driver's license [67]. If a service provider is willing to pay for more expensive verifications—such as those offered through the American Association of Motor Vehicle Administrators (AAMVA) in the US—they may be able to filter out some fraudulent IDs, but even these services have limits, such as not offering comparisons against the photo-of-record [7].

²³ There have already been incidents involving convincing deepfaked video calls of executives [50] and politicians [194].

²⁴ For an overview of Al's ability to solve CAPTCHAs, as well as historical routes to bypass CAPTCHAs, see Appendix B.1.

actions toward goals over extended periods of time, without humans being in the loop or pre-specifying their actions or subgoals [231, 221, 228, 127]. In contrast to existing AI systems—which function mostly as content-generating machines—some of these agentic AIs will be more akin to Internet users: capable of navigating the Web interactively like humans do, but much faster. Some AI developers have predicted that AI agents will be a significant portion of future Internet users. At first, we anticipate that these AIs will be distinguishable by anomaly detection systems that measure atypical Internet behavior on a website. But, over time, AI agents might be able to convincingly mimic the behavioral signals of real human users. With greater dynamism of AI agents will come a larger surface area of risks: For instance, agentic AIs might be able to pursue many dangerous "long con" schemes at once—like creating a number of fraudulent personas to infiltrate open-source communities that manage vital digital infrastructure. Als manage vital digital infrastructure.

AI systems are increasingly scalable

AI systems are not only more capable but also increasingly scalable. Malicious actors can leverage AI tools to execute widespread attacks that they may not have had resources to execute previously. The same factors that lower barriers to beneficial access²⁸ can also enable attacks to increase in scope and frequency.

One factor driving the scalability of AI-powered attacks is cost. AI models are becoming more affordable at every capability level.²⁹ This implies that malicious actors can generate content for a host of deceptive aliases at much lower cost. Thus, with reduced investment required for a successful attack, the frequency of attacks might increase.³⁰ This dynamic is particularly relevant for operations that succeed through persistent and repeated efforts.

Another factor contributing to scalability is the ease of access: The release of highly capable open-weights models makes AI capabilities more available to both well-intentioned and malicious actors. Many open-weights models are available through user-friendly interfaces, decreasing the technical skill required for such uses.³¹ Compared to closed-weights counterparts, open-weights models may increase the relative ease of misuse because they offer less moderation and monitoring of relevant capabilities.³²

²⁵ Upon releasing Meta's latest AI model, CEO Mark Zuckerberg predicted a future of hundreds of millions or even billions of AI agents working online on behalf of small businesses [289].

²⁶There are already instances of AI systems learning to imitate human behavior in complex digital environments [18].

²⁷While this specific scenario is speculative, researchers are evaluating frontier models for signs of dangerous capabilities that could enable such long-con cyberattacks [205, 95, 96]. Social engineering attacks by humans—when they happen over a relatively long timescale—can be very resource-intensive [111, 262]. AI agents may have significant advantages over humans in terms of their levels of patience and ability to multitask across many attacks at once.

²⁸ Some beneficial uses of AI include helping people regain abilities that they have lost, such as using AI to power a voice for someone who has lost the ability to speak fluently [241]. Managing AI risks can be challenging because the same AI applications can be either beneficial or harmful, depending on the scale at which they are used and the intentions behind them.

²⁹ AI systems' abilities have significantly improved in recent years, with a clear relationship between investment and a system's ultimate capabilities, often referred to as "scaling laws" [137, 122, 10, 286]. At the same time that AI systems' abilities have improved, the cost for a given level of ability has generally decreased. See [182] for a visualization illustrating how even relatively cheap AI models from 2024 outperform leading models from 2022.

³⁰ For a framework of different factors that affect Al's impact on influence operations, see [112].

³¹ For instance, some open-weights models are available through desktop applications that can be run on consumer hardware [185].

³² For a recent position paper on the risks and opportunities of open-source generative AI, see [84]. In the context of AI, open-source is often used to indicate releasing the weights and inference code for a model, but not necessarily the full source code by which the model was trained [229]. For this reason, these models are sometimes called "open-weights" rather than open-source.

1.3 Current solutions for countering AI-powered deception need improvement

There are many tools currently used to reduce deceptive and malicious activity online, particularly when the activity is AI-powered. Here, we discuss how the addition of personhood credentials to the toolkit could improve on some tools and complement others, significantly bolstering the foundations of trustworthy interaction. We summarize the discussion in Table 1.

Main deficits
Not robust to highly capable AI
Not inclusive
Not robust to highly capable AI
Not robust to highly capable AI Not privacy preserving
Not scarce

Table 1: Existing tools for countering AI-powered deception and their main deficits.

One approach is **behavioral filters**: distinguishing AI-powered activity based on human abilities or behaviors that are difficult for AI to imitate. For instance, CAPTCHAs historically exploited bots' inability to consistently recognize distorted letters [3]. JavaScript-based browser challenges and anomaly detection systems lean on the fact that—again, historically—bots did not use Web browsers in ways that people do.³³ These filtering methods are becoming less effective as AI systems improve; they are not robust solutions on their own.³⁴

A second approach is to impose **economic barriers** that make it costly to perpetuate AI-powered deception at scale. Some websites introduce paid subscriptions with the explicit intent of distinguishing between people and bots by making it more expensive to use many bot accounts.³⁵ These payments, however, disproportionately affect lower-income users and are an insufficient barrier against particularly profitable forms of malicious activity [110, 245]. Relatedly, some digital services check whether users hold a valid credit card—without necessarily charging a payment to the card—though this approach can be circumvented and also has similar challenges with inclusivity.³⁶

A third approach consists of tools that aim to **detect AI-generated** *content*. Policymakers and technologists have directed investment toward developing "synthetic content transparency" methods in

³³ For instance, the patterns of fraudulent account registrations often differ from those of legitimate users [281]. For a survey of different browser challenges, see [9].

³⁴ In Appendix B.1, we offer a more detailed survey of behavioral filters and their limitations.

³⁵ Elon Musk reported this as one factor in his decision to institute a paid subscription option on Twitter (now X) [224].

³⁶ Credit card requirements can be circumvented via virtual cards [207].

recent years.³⁷ Watermarking, for instance, embeds signals in AI-generated digital content; these signals can help to identify whether content is the output of an AI model. Such tools are useful in many circumstances, but they also have shortcomings: For instance, adversaries may alter AI outputs to try to evade detection, and the tools are not perfectly accurate even if content is undisturbed by an adversary [284, 140].³⁸ More generally, focusing only on whether content is AI-generated does not address other aspects of trustworthy interaction: Consider an army of AI-powered sockpuppet accounts that amplify human-generated content, initially posted from real people's accounts, to push a specific agenda. There is no AI-generated content present—only unobservable decisions about which human-created content to repost—and yet the activity is deceptive.

A fourth approach is to use **document- or appearance-based verification** to confirm that there is a person behind some digital activity. Some verification protocols require that users join a video chat or show evidence that they are in possession of physical identifying documents (for example, they may ask users to send a selfie holding a matching driver's license [244]). At times, these methods go too far toward identifying "who" is behind some activity—collecting more information than required to merely verify that there is "a" person conducting the digital activity. Such approaches are not only potentially intrusive from the user's perspective, but also involve the collection of sensitive personal identifying information, which can introduce important security concerns. Beyond privacy and security issues for users and service providers, such approaches are also not robust to newer AI systems, which are increasingly capable of creating content and avatars that pass these checkpoints.³⁹

Finally, digital service providers often use **identifiers** like emails, phone numbers, and purpose-built **hardware-based authenticators** to try to verify that there is a person taking actions on their services. For example, one-time passwords sent via SMS can offer some indication that the entity behind some digital activity has access to a physical device—though given the rise of tools for acquiring virtual phone numbers and handling complex workflows through these phone numbers, this indication is not particularly strong. There is a more general problem at play: requiring unique phone numbers, email addresses, and even hardware authenticators at signup can reduce some duplicate account creation, but none of these identifiers are scarce enough to establish that their holder is distinct from other users (and thus that the holder could not be creating deceptive identities at scale). Moreover, some of these identifiers facilitate digital tracking [49] and are less private than we aim for with PHCs.

Thus far, we have described trends in AI that could contribute to scalable deception online, as well as reasons why current solutions may be insufficient for countering it. In the remainder of the paper, we articulate a case for personhood credentials as a foundational tool in the toolkit for protecting a trustworthy digital ecosystem, even as AI grows more advanced.

³⁷ For an overview of these methods, see [200]. These techniques were highlighted in President Biden's Executive Order on the Safe, Secure, and Trustworthy Development and Use of Artificial Intelligence [92], which resulted in a report from the National Institute of Standards and Technology (NIST) that provides more detail on these techniques [179]. In Appendix B.2 we offer a fuller discussion of synthetic content transparency techniques, comparing their affordances to those of PHCs.

³⁸ In some open-source model implementations, the watermarking function can merely be removed from the code before running [171].

³⁹ For an example of a particular case in which even a video call did not successfully head off AI-powered deception, see [193]. See also [67, 154]. Beyond using AI to spoof such checks, malicious actors can also purchase photos of people holding up their passports for validation [209].

⁴⁰ For instance, Twilio, a company that provides cloud communications services, has been used to acquire phone numbers used in scams [11].

⁴¹ For more detail on creating fraudulent accounts that require phone number validation, see [248]. Requiring such validation can increase the cost of deception but is not a sufficient solution.

2 Defining personhood credentials

In this section, we outline the design requirements of a personhood credentialing system. Then, we discuss how these requirements balance inherent tensions between preserving user privacy and reducing the possibility of deceptive activity at scale. In Appendix C, we briefly discuss a number of potential implementation approaches—we leave fuller evaluation of such approaches for further research.⁴²

A personhood credential digitally certifies that an issuing entity ("issuer") believes its holder to be a real person who has not received a credential from them previously. The issuer issues a credential through a process we refer to as "enrollment." A third-party digital service ("service provider")⁴³ can request evidence that a user holds a personhood credential as part of some authorization process (like receiving up to a certain number of verified accounts); we refer to this process as "usage." The enrollment and usage processes are illustrated in Figure 3.

Privacy-preserving enrollment and usage of personhood credentials (PHCs)

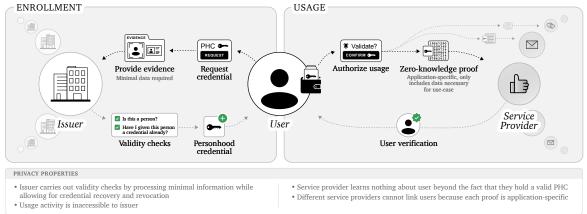


Figure 3: Illustration of enrollment and usage of a personhood credential. **Note:** There are two stages of a personhood credential system. First, the user enrolls. Second, the user can then use the credential with a range of service providers without needing to re-enroll.

When a holder uses their personhood credential, they prove to the service provider that they hold a valid credential without revealing the credential itself (e.g., through a zero-knowledge proof [113]). When necessary, the holder also has the ability to prove that the credential has not yet been used for this particular service (or has been used fewer times than the service's limit). In either case, using a PHC does not reveal any aspect of its holder's legal identity.⁴⁴

⁴² We discuss a range of specific ways to implement credential issuance and usage in Appendix C. One possible implementation option is to build personhood credentials atop the W3C global standards for Verifiable Credentials (VCs), which are a means of expressing a variety of traditional credentials—e.g., driver's licenses, university degrees, passports—on the Web, in a way that is "cryptographically secure, privacy respecting, and machine-verifiable" [240].

⁴³We use "service provider" to mean "relying party," the more common term in security communities.

⁴⁴ PHCs are a type of digital credential, which we define to be an electronic assertion made by an issuer and intended to be used more than once. This definition aligns closely with NIST's definition of digital credential [240]. PHCs could optionally be used in combination with other forms of credentials—like attestations of where one went to college, or where one currently lives—for use-cases that rely on proving further aspects about oneself.

2.1 Foundational requirements of a personhood credentialing system

A PHC system meets the following requirements:

- 1. **Credential limits (1 credential per person per issuer)**: The PHC issuer aims to issue only one credential per person and provides ways to mitigate the impact of transfer or theft of credentials.
 - a. <u>Issuers check one-per-person requirement at enrollment</u>: The issuer has an effective check of whether a person has already received a personhood credential from them.⁴⁵
 - b. Expiry or regular re-authentication: To mitigate the theft or transfer of credentials, there is a periodic process designed to reduce credential use by someone other than the original holder.⁴⁶
- 2. **Unlinkable pseudonymity (privacy)**: PHCs let a user interact with services anonymously through a service-specific pseudonym; the user's digital activity is untraceable by the issuer and unlinkable across service providers, even if service providers and issuers collude.⁴⁷
 - a. Minimal identifying information stored during enrollment: The issuer associates minimum necessary identifying information⁴⁸ between a specific personhood credential and its holder.
 - b. Minimal disclosure during usage: When a user presents evidence of a personhood credential to a service provider, it reveals to the service provider nothing more than "this person holds a valid PHC" or, with the user's authorization, "this person holds a valid PHC not yet used with this service."
 - c. <u>Unlinkability by default</u>:⁴⁹ By default, service providers or issuers cannot trace or link usage activity across uses, even if issuers and service providers collude.⁵⁰ The issuer, by default, learns nothing when a PHC has been used. Service providers do not learn anything when a PHC that has been used with their service is used with another service.

In Appendix C, we elaborate on the assumptions and infrastructure needed to implement a system that meets these requirements.

⁴⁵ Often this check will include some form of evidence that depends on an interaction that occurred in real life, so that AI systems cannot be dispatched to obtain credentials. This evidence *need not* involve a direct in-person interaction between the user and the issuer at the time of enrollment. For example, a passport could be used in an enrollment process, providing evidence that its holder at some point went through the necessary in-person steps to acquire such a document. We offer further discussion of methods for satisfying the credential limit in Appendix C.1.

⁴⁶This regular re-authentication could be achieved through a combination of factors, such as the continued possession of some root document (e.g., a passport), or through tight expiration limits.

⁴⁷We discuss some cryptographic methods that could be used to achieve the requirements outlined here in Appendix C.2.

⁴⁸ The minimum necessary amount of information will vary depending on the types of recovery and revocation procedures the system wishes to afford. Setting tight expiry limits is one way of minimizing the information needed for revocation—see related discussions in a recent comment [20] on the European Digital Identity Wallet Architecture and Framework [88], and in standard ISO/IEC 18013-5 on mobile driving licenses [132].

⁴⁹ By unlinkability, we mean unlinkability via direct use of one's personhood credential. Today, users are sometimes identifiable across websites through factors like their Web browser, even without sharing other identifiers [83]. PHCs are not a direct solution for commonplace tracking of users across the Internet or aggregation of this information via data brokers. However, to the extent that such tracking is facilitated by reusing identifiers across sites (like email addresses or phone numbers), PHCs might make it more difficult to create linkages than previously.

⁵⁰ One risk of government-issued personhood credentials is that the government may be able to compel service providers to turn over records, even if a service provider does not wish to. To be fully privacy-preserving, a personhood credential must be resistant to these forms of subversion.

2.2 Credential limits and the goals of a personhood credentialing ecosystem

A personhood credentialing ecosystem has two core goals, which are sometimes in tension: to reduce scaled deception while also protecting user privacy and civil liberties.

Maintaining a one-per-person per-issuer credential limit helps to balance the ecosystem's goals: attaining reasonable limits on the scale of deceptive activity any person can carry out, while preserving a meaningful degree of choice for users.

We argue that this balance is achieved when there are multiple issuers, each of which limits the number of credentials a person can receive from them. Each person can then obtain a bounded number of credentials—more than one, to counter risks to privacy and civil liberties, but not so many that the credential loses its ability to prevent scaled deception.

Bounded credentials resolve tradeoffs

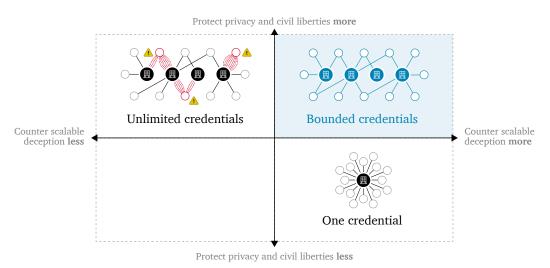


Figure 4: Ecosystem design trade-offs—the argument for bounded credentials.

Note: Open circles represent individuals, and filled circles represent issuers. Top left: Unlimited credentials per person (red). Less effective against deception but better for privacy and civil liberties due to minimal issuer data storage. Bottom right: One credential per person, from one issuer. Effective against deception but risky for privacy and civil liberties. Top right: One credential per person per issuer, with multiple issuers. This balances privacy protection and countering deception.

To illustrate this, we further deconstruct each ecosystem goal into two factors that an ecosystem should try to attain:

Reduce the harms from scaled deception. Limit the scale at which any actor can deceive others, via:

Per-person rate limits. Without per-person limits on some digital activities, harmful actors can pose as multiple people to gain outsized influence on peer-to-peer platforms and in public comment processes.⁵¹

⁵¹ For example, the public comment process for the FCC's 2017 decision on repealing net neutrality protections was compromised by nearly 18 million bot-generated comments [191].

In addition, harmful actors can evade a service's rules—by, for example, creating many new accounts even when caught violating a service's policies.⁵²

Limited incentives for credential transfer. When there are large incentives for theft or sale of credentials, a credential's signal of real personhood erodes over time. Likewise, if it is possible to clandestinely lend one's credential, "authentication farms" may emerge that help bots bypass personhood checks—eventually the credential is not a persistent signal of personhood at all.⁵³

Preserve privacy and civil liberties. Provide users a meaningful choice of issuers and system features that guard their privacy, via:

Minimal information processing and storage. Limiting the amount of personal data processed and stored by credential issuers is essential to preserving user privacy. By minimizing data collection and retention, the ecosystem reduces the risk of misuse or unauthorized exposure of sensitive information.

Checks on power. Preventing the concentration of authority over digital credentials is essential for safeguarding civil liberties, so that no entity can unilaterally dictate terms or exploit personal information stored or processed in the credentialing system.

In Appendix D, we highlight some inherent tensions between these ecosystem design goals through two extremes, which we disfavor: one in which people can acquire an unlimited number of credentials, and one in which people can acquire only one credential from a single issuer. Through these extremes, we illustrate that an ecosystem with bounded credentials may best resolve these tensions. Figure 4 summarizes the argument.

3 Prospective benefits of personhood credentials

We highlight three ways that well-designed PHC systems can reduce the impact of scalable deception online. These benefits are summarized in Figure 5. PHCs can help to:

- Reduce the impact of sockpuppeting: Enable authentic input and engagement from real people at scale, free from deceptive profiles representing people that do not actually exist.
- <u>Mitigate bot attacks</u>: Prevent coordinated attacks of bots circumventing platforms' rules to continue abusing Web services.
- Verify delegation to AI assistants: Signal that an AI assistant is a delegate of a trustworthy person, as opposed to a delegate of a malicious actor.

⁵²The market size of the underground ecosystem for account registration, by one estimate, is roughly 4.8–128.1 million USD per year (as of 2022) [106]. Illicit account creation and the repeat abuse it enables are taxes on digital services and their legitimate users.

⁵³ If personhood credentials can be used with a service without limits (e.g., used to obtain many verified accounts), we expect authentication farms to emerge, in which people verify their PHC with a service and then hand off their account to an AI-powered user. This is comparable to how today's bots can route CAPTCHAs to humans who unblock the bots' path for a fraction of a cent. We discuss these practices in Appendix B.1.

Three key benefits of PHC systems



Figure 5: Summary of three key benefits of PHC systems.

3.1 Reduce the impact of sockpuppeting

Malicious actors lean on "sockpuppetry"—adopting the persona of some hypothetical "person" who could exist but does not—for a range of purposes:⁵⁴ manipulating perception of public political opinion on social media, propping up (or attacking) the reputation of businesses or individuals, and carrying out scams on digital marketplaces [76, 267, 138, 285]. Historically, sockpuppets have been powered either by basic automation or by people in low-wage countries paid by malicious actors to control many misrepresented personas.⁵⁵

Personhood credentials are a mechanism for reducing the viability and impact of sockpuppets. Regardless of whether sockpuppets are AI-powered, services that adopt PHCs can reduce malicious actors' ability to falsely present themselves as multiple individuals on their platform by, for example, imposing account verification limits on each personhood credential. Once a user is verified to be a person, there are many ways a site could incorporate PHCs to facilitate trustworthy discourse, e.g., by boosting or labeling such accounts, or by offering users ways to screen out accounts that are not verified. Scenario 1 outlines one hypothetical use. Such methods are not meant to entirely rid a service of AI-generated content or AI-powered accounts, which can often be benign. But PHCs do offer a range of options for curbing the impact when AI's use is deceptive—particularly at scale.

⁵⁴ Note that many AI-powered personas are benign and beneficial, particularly when transparently disclosed. For instance, an AI-powered chatbot may transparently mimic a science tutor to teach students more effectively. This illustrates a challenging aspect of detecting deceptive uses: An educational institution deploying an AI-powered science tutor is similar in many regards to a bad actor deploying an AI-powered scientist to spread medical falsehoods.

⁵⁵ For one report of deceptive accounts in West Africa purporting to be US-based to influence the US political climate, see [266].

⁵⁶ An account verification limit does not strictly need to be "one verified account per PHC" in order for PHCs to curb malicious activity: A website could decide, for instance, that it wishes for each holder to have at most three verified accounts, which could be unlinked to one another. The important property here is the ability to enforce some finite limit by counting how many times a PHC has been used for a purpose, without more specifically identifying the holder or their particular PHC.

⁵⁷ Stronger distinguishability could facilitate a more trusting online community. On today's Internet, sometimes false accusations are made that a given account is a bot [16]. Knowing another user to in fact be a person could inspire greater trust and ultimately lead to more productive conversations. For instance, some research indicates that using a profile picture with a face in it leads to a greater empathetic response from others in online discussions [157]. PHCs may help achieve a similar effect.

⁵⁸ Some services have established processes for the appropriate disclosure of automated accounts, such as accounts operated with AI or through other computerized systems that do not require human intervention [275].

⁵⁹ While the introduction of PHCs does not entirely solve the problem of AI-powered deception, it dramatically reduces the scale. For example, someone might still hand over their single PHC-verified dating site account to an agentic AI, which they direct to carry out a catfishing scam. Yet the catfisher is less likely to be successful with a single profile as opposed to many potentially targeted profiles used in parallel. Or, someone may post misleading AI-generated content from their single PHC-verified social media account to spread disinformation. But the spreader of disinformation will not be as successful in making their misleading content go viral without amplification from other independent-seeming accounts that they also control.

Motivating Scenario

Authenticating social media accounts.

"Facilitate real human-to-human discussion." That's the goal of Veritalk, a new civic-minded social media platform. How can the platform filter out deceptive profiles generated to mislead—while enabling their users to remain anonymous to both the platform and to other users?

Options like phone number or email verification are no good; it is too easy for bad actors to obtain several and create sockpuppet identities to manipulate the discourse. Same goes for credit card verification, with the advent of virtual cards. Plus, should users really have to pay to be able to have a discussion with a real person rather than a bot? Veritalk aims for inclusivity.

Other platforms verify accounts via government IDs—but Veritalk doesn't like this idea either; too privacy invasive, and malicious actors can use advanced AI to get around these methods with images of fake IDs.

The platform could try to rely on Trust & Safety enforcement—indeed, that may be good in either case, but it is going to get harder and harder to detect AI-powered activity over time.

A member of Veritalk's security team mentions that there is a new state-of-the-art personhood credential available to people in Veritalk's user base. These credentials are exactly what Veritalk has been looking for: a strong assurance that an account represents a single authentic person on their platform, with minimal disclosure of information.

Scenario 1: Authenticating social media accounts.

In recent years, sockpuppets have caused harm across a number of settings. On social media, sockpuppets attempt to manipulate public opinion on a topic [119, 166, 19, 161].⁶⁰ In online dating, bot-powered profiles are used for catfishing, luring users into costly and sometimes dangerous situations [62]. In formal government comment processes, bots overwhelm genuine civic input.⁶²

By reducing the impact of sockpuppets, PHCs might also help to advance a wide range of beneficial experiments in digital governance in the public, private, and nonprofit sectors, though these are speculative and also come with risks. For instance, AI tools might help to digitally gather and analyze public opinions from large segments of the population [98]—PHCs can help ensure that people do not use multiple personas to wield more influence and thereby can increase the perceived legitimacy of such processes. In particular, PHCs can reduce the impact of communities that might otherwise intentionally undermine these processes. Looking ahead, collective input on the design and deployment of AI systems is an

In both cases, if the bad actor is found to be in violation of the service's rules, they can be prevented from returning—a dynamic we discuss further in the following subsection.

⁶¹ Some digital services describe the rationale for their sites' rules in terms of reducing the impact of sockpuppets. X's Community Notes, for instance, states a requirement for accounts to have a unique phone number to "help prevent the creation of large numbers of fake or sock puppet contributor accounts that could be used for inauthentic rating" [274].

⁶⁰ For a broader investigation of sockpuppets in online communities, see [152].

⁶² For example, millions of fake comments were submitted to the Federal Communications Commission's 2017 public request for comment on net neutrality. Researchers have conducted analysis on the methods used to generate such inauthentic comments, which were far more rudimentary than is possible with today's AI [268]. Sockpuppets continue to overwhelm US government public comment processes, and few agencies have taken action to systematically filter out or catch the malicious actors perpetrating these abuses [64, 202]. Recently, a bill aiming to curb foreign and AI-powered abuse in the public comment process for the Bureau of Land Management (the American Voices in Federal Lands Act) was introduced in the US Senate [190].

⁶³ Scholars have argued that new democratic innovations may be especially important in the age of highly capable AI [6, 25]. ⁶⁴ PHCs can help digital governance processes be more robust to claims that AI bots were used to influence a particular

outcome. For more on the ability to cast doubt on outcomes by invoking the specter of AI manipulation, see [225].

⁶⁵ One challenge in soliciting broad input online is that particular communities might flood a poll to try to undermine the poll's intent. For instance, a poll to name a new Mountain Dew flavor produced leading options like "Hitler did nothing wrong" [124], forcing the organizers to take down the naming process. For more, see [219].

emerging special case for soliciting representative opinions, free from sockpuppet manipulation.⁶⁶

3.2 Mitigate bot attacks

Motivating Scenario

Preventing repeated abuse on anonymous platforms.

Likeness Labs is struggling to prevent abuse on their service. Influencers depend on their product to broaden their reach by creating interactive video avatars with customizable personalities and appearances, resembling themselves or fictional personas. These avatars can engage with fans online and integrate with various video messaging services.

However, malicious users are violating usage policies by creating numerous deceptive, high-fidelity likenesses for harmful purposes like political manipulation, financial extortion, or sexual intimidation. Even after account suspension, these bad actors often return with new accounts.

Traditional identity verification or Know Your Customer (KYC) methods are options but can be ineffective, costly, and conflict with Likeness Labs' commitment to user privacy and anonymity—important since some influencers wish to keep their legal identities private.

An alternative is using personhood credentials, allowing Likeness Labs to limit the number of avatars one person can authorize, even across multiple accounts. This enables users to responsibly use digital stand-ins while reducing the feasibility of catfishing and large-scale influence operations.

Although impersonating real individuals other than oneself is against company policy, some users may still do so, depleting their personhood credential's avatar allowance. However, the harm from such impersonations is limited: If a user is caught in problematic impersonation with any of their few avatars, Likeness Labs can suspend their credential.

Scenario 2: Preventing repeated abuse on anonymous platforms.

Malicious actors use botnets—coordinated groups of entities meant to appear as distinct users—to carry out more effective attacks online. Sometimes, overwhelming a Web service with a number of distinct-appearing users is the point of the attack, as in distributed denial-of-service (DDoS) attacks [237]. Other times, the purpose of relying upon many distinct-appearing users is to create many independent chances of success, like circumventing a service's rules and creating new accounts for abuse even after having been caught. As discussed in **Section 1**, without new mitigations, malicious actors powered by AI systems will be more likely to succeed at these attacks.

Personhood credentials can help reduce the impact of bot attacks by enabling service providers to suspend users whose accounts have engaged in abuse. Today, adversaries' tactics—like changing IP addresses and spoofing different browser configurations—often successfully evade the restrictions in place, at least for a period [210]. Some existing software does guard against these attacks, but the developers of cybersecurity software have been locked in a perpetual cat-and-mouse game with attackers. And now, AI changes the dynamic, powering the attackers by reducing the manual work needed to get a new botnet running again after a previous one was caught. By using personhood credentials, service providers can restrict each individual to a limited number of accounts. This means that if someone breaks the rules and their account is suspended, they—identified through their personhood credential used at registration—cannot create new accounts without limit.⁶⁷

⁶⁶ Some AI providers have commissioned research on how to best channel popular input in their tools, including via "Alignment Assemblies" to better align models with democratic will [282, 86, 63, 13].

⁶⁷When establishing restrictions on user activity, websites must balance the level of abuse they are willing to tolerate against the potential loss of legitimate user engagement that such limits might cause [167].

Bot attacks have wide-ranging harms. Genuine users of peer-to-peer digital services are taxed when bots are not kept in check—they face increased friction and are targets of scams.⁶⁸ Digital communities that aim to allow any individual to claim some asset (e.g., free trial memberships, credits for computational costs, products for beta testing), but only once, are often undermined by bot attacks. PHCs can serve as a valuable tool for checking who has already received the benefit, without revealing or tracking any further information.⁶⁹ This rate-limiting can also be applied to mitigate other forms of abuse that are exacerbated by rapid, AI-driven actions, such as ticket scalping, which disadvantages ordinary buyers [250, 187]. Similarly, fraudulent AI-powered requests for benefits from governments or aid-giving organizations—even when successfully defended against—can drain resources and make it more difficult to serve the targets of the aid; screening applicants by PHC verification may help.⁷⁰

Looking ahead, we expect that many forms of AI-powered deceptive behavior that are especially harmful at scale can be *indirectly* mitigated by personhood credentials. In the vignette in Scenario 2, for example, we discuss one particular form of AI-powered deception that is harmful in isolated incidents but especially harmful at scale—impersonation.⁷¹ We dwell on this case because it nicely delineates how PHCs can and cannot directly help when it comes to AI-powered deception online: PHCs by default cannot validate a holder's legal identity, and so a digital likeness service cannot directly check a PHC to confirm that the holder is taking on their own likeness rather than another person's.⁷² However, such activities can be prohibited by digital services' policies, and if a user is ever found to be in violation of the service's policies, PHCs can make the policies more enforceable, by suspending the PHC-linked account and disallowing future signups for that PHC.⁷³ We expect that—anticipating these consequences—fewer attackers will attempt such unauthorized impersonations.

3.3 Allow verified delegation to AI agents

AI systems are increasingly capable of acting autonomously [141, 264, 95, 277]. While enabling many beneficial use cases,⁷⁴ this autonomy ("agenticness") also facilitates a new form of deception: bad actors can deploy AI systems that, instead of pretending to be a person, accurately present as AI agents but pretend to act on behalf of a user who does not exist. This strategy exploits the current lack of norms around disclosure of agentic AIs, including a lack of norms around disclosing the identities of the people controlling them (often called their "principals").⁷⁵

⁶⁸ For instance, a website may have particular processes—phone number verification, rate limits—that are primarily meant to limit abuse by repeat bad actors but that inadvertently create difficulties for benign users.

⁶⁹ A prosaic example of such fraud is free trials for digital services, which are often serially abused [77].

⁷⁰ One risk is that screening via PHCs may reduce fraud but also make it more difficult for legitimate recipients to receive aid: Administrative overhead to apply for benefits tends to result in many people not receiving benefits that they are properly eligible for [120, 226].

⁷¹ By impersonation, we mean taking on the likeness of someone in particular, as opposed to taking on a sockpuppet persona—a plausible human that could exist but does not.

⁷² Malicious actors may have incentives to impersonate a real person to carry out social engineering attacks or to set up a profit-making scheme (e.g., if the target of impersonation is a celebrity).

⁷³ Beyond violating the policies of a digital service, using AI for impersonation might also sometimes involve breaking a law. Courts have sometimes ruled in favor of a "right to publicity" and "right to privacy," which using a person's likeness without their permission may violate, depending on the circumstances [269, 101].

⁷⁴ For a description of some economically useful properties of AI agents, see [231]. Already, some humans are deferring to AI-powered solutions for navigating dating apps on their behalf [230].

⁷⁵ In some jurisdictions there is already regulation governing bot disclosure [143]. In 2023, legislation was introduced to codify bot-disclosure federally, though it has not become law [189, 258]. Note that these bills outline requirements around disclosure of AI usage, i.e., requiring that people disclose when a bot was employed to take an action or when a piece of content was generated by an AI model. The outlined disclosure requirements stop there—they do not include disclosures related to the

Motivating Scenario

Accountable delegation to AI assistants on anonymous platforms.

A new marketplace website, TradeHub, enables peer-to-peer transactions of goods and services in a trustworthy but privacy-preserving way—like the classified advertisements you used to find in the newspaper, but for the digital age.

Users' AI assistants could post listings, respond to posts, negotiate deals, and schedule meet-ups. This meant that users could maintain an active presence on the platform without constantly monitoring their accounts.

All was going well at TradeHub, but as AI became more advanced, the site started observing a new category of fraud: AIs that seemed to be using the site on behalf of a supposed human user, who the site was not sure actually existed.

Typical security measures like KYC checks were too cumbersome and unappealing for TradeHub's diverse users. They needed a solution that ensured trust and accountability without sacrificing privacy or convenience.

Enter personhood credentials. These allowed TradeHub to discern which AI assistants were backed by real people, without exposing personal details to the platform or other users.

TradeHub users loved this addition. With personhood credentials, they could filter out unreliable AI assistants and prove their own AI helper was trustworthy, making anonymous peer-to-peer exchanges smoother and more reliable than ever.

Scenario 3: Accountable delegation to AI assistants on anonymous platforms.

Personhood credentials could offer a way to verify that AI agents are acting as delegates of real people, signaling credible supervision without revealing the principal's legal identity. This feature could be useful in a range of settings where users wish to rely upon AI assistants; Scenario 3 outlines one such use case. Should a principal fail to address harms caused by their PHC-verified AI, they risk suspension from a service. Suspension implies that they lose their ability to verify delegates for some time period, reducing their capacity to perpetrate future harms.

Note how, in this case, PHCs create a form of accountability for AI agents without demanding sensitive information from principals.⁷⁶ This suspension mechanism may effectively signal which AI agents have trustworthy principals, even if verifying AI agents through principal-linked PHCs is voluntary. Agents that remain unverified might be perceived as having reasons for not undergoing verification.⁷⁷ Moreover, some malicious activities involving autonomous AI agents may rely on hiding the fact that multiple agents are controlled by the same individual.⁷⁸ PHCs can help address this issue by creating a framework

identity of the person or organization enlisting the AI.

⁷⁶ A large and growing body of legal scholarship explores when and how humans should be accountable for the harms caused by AI agents [87, 54, 55, 52, 56, 156, 74, 144]. Many—though by no means all [238, 100]—of these proposals involve holding humans liable for some harms caused most directly by those agents. It is beyond the scope of this article to recommend any of these approaches over others or to explore the finer points of how such theories could work. We note, however, that these theories rely on the ability of someone harmed by an AI agent to sue the principal of that agent, which in turn depends on the principal being *identifiable*—which a PHC alone does not achieve. Even without directly identifying the principal, however, PHCs can still shift the benefits of AI agent usage to be more positive.

⁷⁷ This mechanism aligns with the "unraveling result" of voluntary disclosure models in economic theory [172, 173]. In scenarios where parties hold verifiable private information—such as whether an AI agent operates under a trustworthy principal—even if revealing this information is optional, agents associated with trustworthy principals have a strong incentive to disclose it. Consequently, a lack of disclosure becomes informative: in equilibrium, agents that do not disclose are effectively signaling that they do not have trustworthy principals.

⁷⁸ One example of how a person might enlist unverified agents to manipulate a market: A potential buyer on an auction site might create multiple personas to submit lowball offers, in hopes of baiting the seller into an artificially low assessment of

that links multiple AI agents to a single principal without revealing the principal's specific identity. By doing so, PHCs might make it more difficult for bad actors to conceal their network of AI agents, thereby reducing the potential for abuse that stems from undisclosed common principals.

Sometimes, a website—or a third-party user interacting with the agent—may need to verify the Al's specific principal, not merely that it is backed by *some* principal. Ultimately, a fuller framework for verifying AI agents [43, 44] and their principals will likely be necessary.

4 Prospective challenges for personhood credentials

In this section, we discuss four areas in which PHCs' impacts must be carefully managed:

- Equitable access: How might PHCs impact access to digital services?
- Free expression: How might PHCs impact whether people feel safe and confident engaging across digital services?
- Checks on power: How might PHCs change the power dynamics of digital services, in both substance and perception?
- Robustness to attack and error: How might PHCs be vulnerable to mistakes and intentional subversion by different actors in the PHC ecosystem?

We highlight these areas as most salient among our discussions with security researchers, civil liberties groups, and builders of digital identity systems, and acknowledge that they cover only a small subset of all possible unintended consequences. This section is intended to be taken only as a preliminary outline. The risks discussed in this section can be managed through the pursuit of ideas discussed in **Section 5**.

4.1 Equitable access

Benefits from personhood credentials must be weighed against potential impacts on the accessibility of digital services. A need for frequent authentication of one's PHC could contribute to friction and frustration for users, ⁷⁹ particularly those less comfortable with technology generally (like older adults, who are common targets of AI-powered deception [246]). ⁸⁰ If friction from PHCs is too high, some users may neglect to use them and thus lose out on the range of benefits they afford.

their item's value. In a similar example, a potential seller might enlist synthetic personas to create the appearance of higher demand for their item, inducing a legitimate buyer to pay a higher price and to act with more urgency. In each case, the ability to recognize these entities as linked to the same identity—even without knowing any further details of the identity—would help to head off these risks. For an estimation of the prevalence of these practices, known as "shill-bidding," in online markets, see [51]. For a discussion of market designs that are robust to these practices, see [145].

⁷⁹ Although much of this paper focuses on the importance of preserving the Internet's usability for *people*, many automated software scripts are also important for a well-functioning Internet. Such scripts could be inadvertently affected by the widespread introduction of PHCs; websites may want to consider approaches that offer benign bots a way to proactively distinguish themselves. One example of how some websites communicate their preferences for bot interaction is the use of robots.txt files, which have existed in some form since the 1990s as an opt-in approach for contending with the impacts of bots navigating the World Wide Web [149, 57].

⁸⁰ It is well recognized in commercial contexts that the friction of authentication is typically a negative experience for users [109, 114]. Depending on how authenticating with a PHC works, "Sign-in with PHC" options may ease some of this friction.

Relatedly, if service providers use personhood credentials to limit access to particular services, some groups of individuals without PHCs may be systematically excluded. It is thus important that PHC issuers design their systems with equitable access in mind, and that authorities guiding PHC ecosystem development are careful to ensure that all groups have some form of access to vital services. Depending on the details of a PHC's implementation, some groups may have access before others: It is important that this not contribute to the economic exclusion of slow-adopting groups in the meantime. One positive possibility is that PHCs may help to unlock digital services for many who do not have access to forms of digital identity and are thus currently excluded from many basic digital services. When assessing PHCs' relative impact on accessibility, it is also important to consider that as AI systems become increasingly indistinguishable from real people online, individuals and service providers may default to mistrusting digital strangers, which could likewise disproportionately affect some groups.

To minimize their impacts on accessibility of digital services, PHCs must be robust to changes in a person's circumstances. For instance, PHCs must be robust to changes in citizenship or residence. There must be fallback processes for PHCs that are not too taxing on the user, and which are extensively trialed to identify edge-cases. Further, PHC providers should recognize that a fear of exclusion due to life changes may be psychologically taxing, even if there are appropriate measures to regain access.⁸⁵

4.2 Free expression

It is important to consider how PHCs could reduce people's willingness to speak freely and to dissent in digital spaces. Individuals may fear that their online speech will be linked to their offline identity, even though PHCs do not convey one's identity.⁸⁶ Both perceived and actual anonymity are important for expressing views without fear of retribution, particularly for marginalized groups or for those living under authoritarian regimes with extensive surveillance [206, 153].⁸⁷

PHCs also change the dynamics of credibility online, possibly in unpredictable ways that could affect individuals' willingness to speak openly. Individuals who choose not to verify their personhood through PHCs might find their contributions discounted or even labeled as disinformation.⁸⁸ On the other hand, statements that might previously have been dismissed due to doubts about their origin could earn credibility when accompanied by PHC verification. In some cases, hostile governments might be able to

⁸¹ One analogous precedent may be found in the United States's Blueprint for an AI Bill of Rights, which states that certain vital services should provide a way for users to opt out of AI interaction; see "Human Alternatives, Consideration, and Fallback" in [188].

⁸² For example, if PHCs are primarily backed by government IDs, there will likely be geographic disparity in access to PHCs initially. Around 850 million people worldwide do not currently have an official identification document [272]. Refugee populations, in particular, may struggle to obtain or maintain a PHC. This issue must be investigated, perhaps under the purview of existing programs studying refugee documentation, e.g., within the United Nations [257].

⁸³ PHCs will be less likely to contribute to economic exclusion if they are an opt-in mechanism for additional components of a service, rather than mandatory for the baseline service itself.

⁸⁴It is estimated that 3.3 billion people worldwide do not currently have access to an official digital identity for online transactions—PHC systems could offer an easier-to-access alternative [272].

⁸⁵ Processes for changing one's name, gender, and other details can be very cumbersome. For instance, in the United States, a person must contact a host of government agencies when changing their name [261].

⁸⁶ PHCs therefore respond to key criticisms of some platforms' "real name" policies [30, 279].

⁸⁷ Anonymity has been referenced as an essential part of public discourse. The United States Supreme Court stated in *McIntyre v. Ohio Elections Commission* that anonymous speech "is not a pernicious, fraudulent practice, but an honorable tradition of advocacy and of dissent" [164]. See also [148].

⁸⁸ This dynamic has some similarity to the concept of the liar's dividend [225]. Though PHCs are implemented with the intent to mitigate online disinformation, they might have the side-effect of making non-PHC verified speech easier to write off as disinformation, even if it is in fact authentic.

suppress authentic and credible dissent by limiting access to PHCs within their borders.⁸⁹

While PHCs preserve user privacy via unlinkable pseudonymity, they are not a remedy for pervasive surveillance practices like tracking and profiling used throughout the Internet today. Although PHCs prevent linking the *credential* across services, users should understand that their other online activities can still be tracked and potentially de-anonymized through existing methods [83]. It is important to recognize that if PHCs are challenging or inconvenient for users to use, companies may be inclined to direct them toward alternatives that are easier to use but provide less privacy. Consequently, some users might opt for these more user-friendly, yet less private, solutions. Emergent social dynamics may also threaten the anonymous nature of PHCs—for instance, if online abuse continues to be persistent even among individuals with PHCs, people may feel socially pressured to reveal auxiliary personal information to demonstrate that they are trustworthy.

There are many nascent legal and ethical discussions centered on the rise of highly capable AI and freedom of expression. It will be crucial to understand how PHCs fit into the legal and ethical frameworks that emerge. For example, one important legal question is when individuals need to disclose their use of AI [143, 189, 258]. Service providers might wish to allow AI delegation on their platform only with proper transparency to other users, or they might wish to disallow AI delegation on their service altogether. PHCs could provide one method for disclosing these new forms of expression, as well as a means of enforcing the related rules that do emerge.

4.3 Checks on power

One significant challenge for a PHC ecosystem is how it may concentrate power in a small number of institutions—especially PHC issuers, but also large service providers whose decisions around PHC use will have large repercussions for the ecosystem. ⁹³ Such challenges may be particularly acute in cases where the PHC issuer is not a democratically elected government with checks on its powers and accountability to its constituents, or where influential service providers are operating with little regulatory oversight.

User information is one potential source of power for an issuer: Issuers should give strong assurances about what information they hold and for what reasons, favoring an "as little as possible" approach and confirming compliance via mechanisms like audits. Transparency over technical aspects of the system (such as open-sourcing certain software components⁹⁴) may also reduce the risk that data are misappropriated⁹⁵ and the perception of such risks. Possible mitigations may include granting users a

⁸⁹ There are a number of ways that a hostile government might try to stop dissidents from getting PHCs. For instance, if a PHC system relies upon proprietary hardware, a hostile government might be able to restrict the hardware from its country, or otherwise control access to the hardware so that the government gains control over users' PHCs [38]. Governments could also apply these prohibitions more narrowly to certain groups.

⁹⁰ Such practices are widely documented; for one comprehensive treatment, see [288].

⁹¹ This is already seen with some examples of digital identity protocols and cryptocurrency wallets [146].

⁹² For example, a dating application may wish to disallow users delegating any conversations to an AI agent if they have a "verified person" account. Nonetheless, the service provider may have a difficult time determining that the user has illicitly delegated to an AI agent, if this delegation occurs after signing up using one's PHC. If services wish to curtail delegation—particularly when users have incentives to illicitly delegate—services will need methods for detecting signals of AI-powered behavior (e.g., if a user is on the service at all hours of the day) and for enforcing rules against it thereafter.

⁹³ Appendix D contains further discussion about how different design choices at the ecosystem level may produce power asymmetries between users and issuers, as well as between users and service providers.

⁹⁴ For example, Signal's server [236] and protocol libraries [235] are open-source, and their technical specifications are publicly available (see, e.g., [160]).

⁹⁵ Even companies with large financial consequences at stake have suffered breaches of sensitive user data [102].

"right to be forgotten" [270, 136] by a PHC system—for instance, if they have lost trust in the issuer.

The range of decisions under issuers' and service providers' purview may also create imbalances of power relative to users. An issuer may be able to choose which use-cases and service providers to support; 96 similarly, service providers may be able to choose which issuers of PHCs to accept as legitimate. These choices might be made for reasons related to private incentives, rather than the well-being of users or an overall community. Furthermore, service providers may have opportunities for illegitimate exercise of their power through their increased ability to enforce controversial policies. 97

Proper democratic oversight, accountability, and transparency mechanisms must be in place to check the power of issuers and service providers, whether the issuers are governments or nongovernmental entities. More broadly, the incentives of the issuer must be carefully considered—for instance, some existing PHC systems are associated with particular cryptocurrency tokens, which may introduce complex financial relationships between the issuer and credential holders.

4.4 Robustness to attack and error

A PHC system, like any digital system, is vulnerable to attacks and exploits by multiple actors—most notably, subversion by the issuer itself, by service providers, and by users with malicious intent (for instance, those in a network of cybercriminals). Many cybersecurity best practices will apply to these systems—for instance, defenses against denial-of-service attacks that halt the issuance of new personhood credentials, and means to stop attackers from gaining access to sensitive records. ¹⁰⁰

One important threat to consider is whether a PHC system is robust to subversion by its issuer. For instance, issuers may have asymmetric knowledge about vulnerabilities in their system, which allow them to issue fraudulent credentials for self-interested uses. One must also consider the possibility that a PHC issuer becomes less trustworthy over time—whether due to personnel changes or other factors. PHC systems may wish to commit to practices like third-party audit compliance to demonstrate their trustworthiness and ensure that users can discover if trust is no longer merited.

Another important threat to consider is how malicious service providers may use PHCs as cover to surreptitiously collect information about users. For instance, even though PHCs do not disclose any aspects of a user's legal identity, a service provider might be able to trick users into uploading pictures of their government IDs to (purportedly) re-authenticate their PHC.¹⁰¹

⁹⁶ In some identity systems, like India's Aadhaar, private-sector service providers need to apply to the government to be authorized to identify people via their Aadhaar number [255].

⁹⁷ Just as services can use PHCs to stop bad actors from circumventing suspensions via new accounts, such enforcement could also impact benign users who have accidentally violated a service's rules. For services that use PHCs to enforce their rules, it is particularly important that users can discover and understand services' rules ahead of time.

⁹⁸ In some instances, governments have restricted data collection by issuers of nongovernmental proof-of-personhood systems like Worldcoin, claiming that data collection practices violated local laws [184] (some restrictions have been subsequently lifted [215]).

⁹⁹ Historically, one motivation for proof-of-personhood systems—closely related to PHCs—has been to aim for equitable representation when voting in digital communities, like those associated with blockchain communities. Prominent examples of proof-of-personhood protocols with an associated currency include Idena, Worldcoin, Proof of Humanity, and BrightID [72]. A fuller evaluation of existing proof-of-personhood systems, and whether the credentials they issue meet the PHC requirements (credential limits and pseudonymous unlinkability), is beyond the scope of this paper.

¹⁰⁰ The potential harm of data extraction attacks can be reduced through strict adherence to the requirement, detailed in Appendix C.2, that PHC issuers store only the minimum necessary information.

¹⁰¹ Standard authorization and re-authentication processes across services might reduce the ability for service providers to

A final threat we consider is how users may try to subvert a PHC system to obtain multiple credentials, either at small scales or through coordinated attempts at cybercrime. These users could attempt to deceive the enrollment process, or could rely upon theft or illicit purchase of other people's credentials. ¹⁰²

Beyond the attacks just described, a PHC system should also be robust to user error. For instance, PHC systems can offer methods to ensure that credentials can be recovered or revoked in case they are lost or stolen. A user should not be permanently locked out from vital digital services if they merely misplace their smartphone or other means of accessing their PHC; it is incumbent upon system designers to make PHCs error-tolerant for benign users.

These threat models and considerations are by no means comprehensive; we encourage research that classifies and assesses the many attacks that PHC systems may face, which will help inform the design of more robust systems.

5 Next steps for consideration

Governments, technologists, and standards bodies—with frequent feedback from the public—can take steps to manage the risk of AI-powered scalable deception. Here, we outline two broad categories of ideas: **adapting** existing digital systems to prepare for the impacts of highly capable AI, and **prioritizing** the development of personhood credentials as one specific countermeasure.

Next steps					
Governments, technologists, and standards bodies could, in close consultation with the public, consider the following ideas to address the threat of Al-powered scalable deception online:					
Adapt	A1	Reexamine standards for remote identity verification and authentication.			
existing digital systems A2	A2	Study the impact and prevalence of deceptive accounts on major communications platforms.			
	A3	Establish norms and standards to govern agentic AI users of the Internet.			
Prioritize personhood credentials	P1	Invest in the development and piloting of personhood credentialing systems.			
	P2	Encourage adoption of personhood credentials.			

Table 2: Summary of next steps to consider.

surreptitiously collect information through non-standard means.

¹⁰² One analog for these concerns is the market for resale of fraudulent accounts on social media sites [247].

¹⁰³ NIST offers recommendations on the management of cryptographic keys, which may be applicable to the design of recovery and revocation processes for PHCs [134]. The easier that it is for a holder to recover or revoke a compromised PHC, the less incentive an attacker will have to try to compromise a credential, as they should not expect to control the compromised PHC for long. On the other hand, if it is too easy to revoke one's PHC and be reissued a new one, this may reduce the issuer's ability to enforce its credential limit requirement of one per person per issuer. We discuss this challenge further in Appendix C.

5.1 Adapt existing digital systems

As we have discussed throughout the paper, many features of the Internet and digital infrastructure may struggle to contend with malicious actors whose deceptive schemes are powered by increasingly capable AI. It is important to assess the Internet's readiness for these AI systems and to adapt accordingly.

A1. Reexamine standards for remote identity verification and authentication.

The emergence of highly capable AI technologies may render certain security assumptions inaccurate. Such assumptions may be common to a range of standards that governments and industry actors rely on. Key standards that may require reevaluation include NIST's Identity Assurance Levels (US), ¹⁰⁴ the European Union Agency for Cybersecurity's (ENISA's) practices for remote ID proofing, ¹⁰⁵ and the international standard ISO/IEC 29115 regarding entity authentication [131].

Anti-spoofing recommendations are one component of standards that may now be outdated: Previously, asking a subject to replicate a certain pose on a live video call was reasonable evidence that they were a real person rather than a prerecorded video of one. But malicious actors can increasingly use AI to mimic actions like these with lifelike avatars [147, 33, 276, 259, 123].

The relative strength of different forms of evidence may also warrant reevaluation based on Al's improved abilities. For instance, authentication methods once considered effective in reducing fraud—such as voice-based authentication—are now increasingly vulnerable to AI-driven attacks. Advanced voice synthesis technologies can replicate an individual's tone and speech patterns using minimal data, turning voice authentication into a potential avenue for account takeovers [196]. In some domains (e.g., images and videos), the reliability of a piece of evidence might vary according to whether there is proof that the media was captured directly from a camera. More generally, reliability assessments of different forms of evidence may be improved by directly factoring in how AI might be used to subvert the evidence. 108

A2. Study the impact and prevalence of deceptive accounts on major communications platforms.

Currently, there are no standardized, official methodologies or tools to reliably measure the prevalence of AI-powered accounts across various digital platforms. Developing such tools could enable platforms and governments to perform robust comparative analysis and to determine the effectiveness of different mitigation strategies. Encouraging independent auditing firms to assess the presence of AI-generated

 $^{^{104}}$ NIST's Digital Identity Guidelines are actively being revised [134, 181].

¹⁰⁵ ENISA's 2024 report includes a full section on deepfake-powered attack vectors, though it also suggests certain indicators of deepfakes that we do not expect to persist—for instance, certain characteristics about their bitrate or frames per second [91]. ¹⁰⁶ For instance, NIST's guidance on Identity Verification suggests: "In order to confirm the video stream is live and not pre-recorded, the operator may direct the applicant to move their head in specific ways, raise or lower eyes, or ask the applicant questions requiring response during the live capture video" [180].

¹⁰⁷ Photos taken with a C2PA-signing camera, like the smartphone application Truepic, can attest to being untampered [254]. ¹⁰⁸ For instance, NIST considers biometric verification to be its highest strength of evidence, even if done via remote identity proofing [180]. To remain confident in the strength of this evidence, standards bodies should consider explicitly stress-testing any biometric algorithms for robustness to deepfake attacks: As NIST notes, "The biometric False Match Rate (FMR)...does not account for spoofing attacks" [117].

¹⁰⁹ Without a standard approach, some have speculated that individual platforms have an incentive to permit some AI-powered accounts to drive up the total number of users and engagement, thus allowing the platforms to charge more money for advertising. For instance, ahead of Elon Musk's acquisition of Twitter, he accused the site of engaging in these practices [69]. The company that Musk commissioned to conduct a study of Twitter's bot activity—Cyabra—has claimed to have benchmarks across social media companies, but to our knowledge has not made these data or methodologies available [79].

content and accounts might promote transparent and unbiased analyses across communications platforms. Moreover, restoring or enhancing researchers' access to platform data can offer insights into the proliferation and trends of AI-driven accounts. 110

Governments might explore several approaches to increase the availability of this information. For example, they could establish incentive programs for entities that share relevant information or that fund research grants focused on this topic. Clearly, the scope of such government activities would need to operate within existing legal protections of speech on private platforms, ¹¹¹ as well as user expectations of privacy.

A3. Establish norms and standards to govern agentic AI users on the Internet.

We anticipate that, in the near future, AI agents will constitute a much larger portion of activity on the Internet than they do now [53]. Appropriate trust in these agents can be facilitated by systems for confirming agents' identities, properties, and claims (e.g., about which data and applications they have permission to access) [43].

As AI agents begin to interact with digital services in new ways, it is important to establish safeguards and guidelines around permissible interactions. Linking AI agents to personhood credentials could be a starting point, but a more comprehensive agent certification framework may be necessary for certain applications. Many organizations are already working on best practices for issuing, authenticating, and presenting digital credentials—the EU Digital Identity Wallet Consortium [88] is a prominent example—and these initiatives can more directly prepare for highly capable AI.

5.2 Prioritize personhood credentials

We encourage governments, technologists, and standards bodies to prioritize the development, piloting, and adoption of personhood credentials (PHCs) to reduce AI-powered scalable deception.

P1. Invest in the development and piloting of personhood credentialing systems.

To be useful, PHC systems must attain a meaningful scale.¹¹⁴ One way to build trust and spur adoption is to build PHC systems incrementally atop existing government credentials, such as driver's licenses and passports.¹¹⁵ We expect that public-private partnerships will be important for moving nimbly: A number

¹¹⁰ Some social media companies have phased out these forms of researcher access programs in recent years. In 2023, Twitter (now X) ended free access to APIs that had enabled research such as estimates of bot prevalence on the website [39]. More generally, some social media companies have taken action to restrict research programs studying dynamics on the sites [28].

¹¹¹ In the United States, Section 230 of the Communications Decency Act is especially relevant when considering platforms' responsibility for speech on their sites [24].

¹¹² For suggestions of how to govern increasingly agentic AI systems, see [231].

¹¹³ The development of HTTPS infrastructure may serve as an informative case study. The HTTPS certificate ecosystem uses a series of trusted institutional checks and public-key cryptography to allow users to trust that they have been routed to the correct website [81].

¹¹⁴There are some public estimates of the scale of proof-of-personhood systems, which are closely related to PHC systems. The largest of these, Worldcoin, announced in July 2024 that it had reached 6 million signups [273]. Idena, another proof-of-personhood system, appears to have had on the order of several thousand members as of April 2022 [129, 192]. The proof-of-personhood system Proof of Humanity had roughly 17,000 members as of 2022 [170].

¹¹⁵ Such work might fit naturally into an ongoing initiative led by NIST in the US related to digital driver's licenses and the management of digital identity on mobile devices generally; see [168] for write-up and [177] for project status. Certain types of digital driver's licenses have different privacy properties, which may afford different functionality when used as a basis for a

of technologists have already developed proofs-of-concept for converting government identification documents into privacy-preserving credentials¹¹⁶ and may be amenable to expanding these pilots. Notably, some regional governments, such as the government of British Columbia, are already issuing similar credentials [34, 165]. If many governments or multinational organizations pursue PHCs, a supranational entity could help to coordinate PHC systems internationally.¹¹⁷

Accelerating research and development of PHC systems—particularly their privacy, inclusivity, and fraudresistance—is essential. The defenses of PHC systems can be assessed through red-teaming exercises that identify vulnerabilities and point toward countermeasures [283]. Inclusivity research might benefit from exploring the adoption barriers faced by various demographic groups.¹¹⁸

There are many open questions, particularly legal ones, surrounding personhood credentials:

- How should personhood credentials relate to existing identity theft and protection laws?
- In scenarios where service providers currently cannot require government ID, should similar restrictions apply to government-issued personhood credentials?¹¹⁹
- Are there specific use-cases that warrant granting individuals the right to maintain multiple verified pseudonymous accounts, as opposed to allowing service providers to mandate a single verified account per person?
- How can users confirm that privacy commitments from an issuer and service provider are being upheld, without relevant cryptographic expertise?

Research that explores these issues will help stakeholders evaluate implementation trade-offs. Moreover, people can then engage in substantive dialogue that enhances PHCs through broad-based deliberation.

P2. Encourage adoption of personhood credentials.

Encouraging service providers to accept personhood credentials (PHCs) as alternatives to traditional ID verification—especially in scenarios where legal identity is not strictly necessary—could substantially boost PHC adoption and enhance privacy. Governments might also define situations where individuals have a right to interact with only other real people, highlighting PHCs as one tool for achieving this.

To foster PHC demand, governments can promote their use in digital civic activities, like public comment processes. These uses not only normalize PHCs but also showcase their practical benefits, fostering broader public acceptance and trust. Additionally, governments can consider offering subsidies or incentives targeting user enrollment and inclusive outreach efforts.

PHC. For criticism of the mDL standard (ISO/IEC 18013-5) for digital IDs, see [242]. A digital ID's validity period is one factor that might affect the potential ability to triangulate users' activity within a certain time window, and therefore affect the ability to satisfy the PHC requirement of pseudonymous unlinkability [20].

¹¹⁶ For two examples of protocols that convert an existing identity document to a privacy-preserving credential, see [287, 12]. More broadly, the US government has funded infrastructure solutions that demonstrate support for privacy-preserving digital credentials [227].

¹¹⁷ Such an entity could focus on cross-border interoperability and standards, akin to how the UN's International Civil Aviation Organization manages standards for biometric passports [130].

¹¹⁸ For example, policy initiatives related to broadband access may serve as an instructive precedent—in that case, bridging the digital divide required focusing special attention on rural communities, people with disabilities, and underserved groups [260, 139].

¹¹⁹ As an example of restrictions on existing government ID programs, India's Supreme Court ruled that its banks may not demand Aadhaar details as a mandatory process for opening a bank account [159].

¹²⁰ For an illustration of the risks of sharing government IDs for verification purposes, see [66].

To facilitate PHC adoption, it is valuable to improve general digital accessibility and literacy. Pre-installing digital credential management applications on smartphones and highlighting them as key features can aid users. ¹²¹ Encouraging broader use of tools like password managers can prepare individuals to secure their personhood credentials effectively. ¹²² As users become more comfortable with these technologies, they are more likely to enroll in PHC systems and to use PHCs successfully.

¹²¹ The ability for iOS and Android devices to store digital identity cards in their Wallet applications is one possible start for

¹²² Greater understanding of concepts like encryption and zero-knowledge proofs may also help people understand the ways their personhood credentials are secured—easing some anxiety about whether PHCs' use is truly private.

Acknowledgments

This research project and its resultant paper—focused on the challenges that AI could pose for trustworthy digital interaction, and on prospective solutions to these—are the efforts of a large multidisciplinary collaboration, across researchers at many different organizations, with many different points of view and equities. For instance, coauthors hold a range of views on the relative merits of different "roots of trust" for a PHC—whether to build on government IDs, Webs of Trust, and/or biometrics. These disagreements center on complex empirical and normative questions around the robustness, inclusivity, and privacy risks of different methods. Likewise, coauthors hold a range of views regarding whether an ideal PHC system should offer enrollment to people worldwide or focus more narrowly on a certain geography. One objective of this paper is to catalyze public discussion and debate about these topics. This paper should not be read as implying the full point of view of any particular coauthor or of the organization to which they are affiliated. Further, this paper should not be read as advocating for or against any particular PHC system or digital credential.

We are grateful to a number of groups for feedback and input that helped to shape the ideas and presentation of this paper. We would like to thank Sandro Herbig of Tools for Humanity¹²³ for his contributions to these ideas, including detailed explanations of proof-of-personhood systems' functionality and common points of feedback and confusion related to these systems. We would like to thank Eden Beck and Erol Can Akbaba for copyediting and typesetting, respectively, and for generally sharpening the presentation of the paper's ideas.

We would like to thank the AI and Media Integrity Steering Committee of the Partnership on AI and the Berkman Klein Center for Internet & Society at Harvard University for hosting us for feedback sessions related to these topics. We would also like to thank a number of individuals for feedback and conversations: Scott Aaronson, Mada Aflak, Dan Alessandro, shirin anlen, Rahul Arora, James Aung, Boaz Barak, Katie Benjamin, Jake Brill, Miles Brundage, Jon Callas, Rosie Campbell, Alan Chan, Melissa Chase, Peter Cihon, Paul Crowley, Karen Easterbrook, Tyna Eloundou, Emiliano Falcon-Morano, Jason Fedor, Claudia Fischer, Tim Fist, Lydia Gorham, Mark Gray, Margaret Hu, John Jordan, Gabe Kaptchuk, David Karger, Pedram Keyani, David Kim, Grace Kwak Danciu, Kim Laine, Michael Lampe, Jaron Lanier, Teddy Lee, Traci Lee, Brenda Leong, Leon Maksin, Timothy Marple, Stephen McAleer, Reed McGinley-Stempel, Pamela Mishkin, Ben Newhouse, Cullen O'Keefe, Christian Paquin, Joel Parish, Tobias Peyerl, Brendan Quinn, Ankur Rastogi, Francis Real, Elizabeth M. Renieris, David Robinson, Ben Rossen, Girish Sastry, Bruce Schneier, Eric Scouten, Gaurav Sett, Yonadav Shavit, Dane Sherburn, Adam Shostack, Allison Stanger, Jay Stanley, Esther Tetruashvily, Değer Turan, Raquel Vázquez Llorente, Chelsea Voss, Becky Waite, E. Glen Weyl, Rebecca Williams, Wenda Zhou, and Jonathan Zittrain.

¹²³ Tools for Humanity is the technology company that led the initial development of Worldcoin / World ID—an instance of a proof-of-personhood system—and works to support the World ID project's technology protocols. One of Tools for Humanity's co-founders, Sam Altman, is also co-founder and CEO of OpenAI. This paper is not an endorsement of World ID or any specific personhood credential or proof-of-personhood system.

References

- [1] ABC4Trust, n.d. Attribute-based Credentials for Trust. URL https://abc4trust.eu/. Accessed: 2024-07-10.
- [2] M. Acquah, N. Chen, J.-S. Pan, H.-M. Yang, and B. Yan, 2020. Securing fingerprint template using blockchain and distributed storage system. *Symmetry*, 12(6):951.
- [3] L. Ahn, B. Maurer, C. McMillen, D. Abraham, and M. Blum, 2008. reCAPTCHA: Human-based character recognition via web security measures. *Science*, 321(5895):1465–1468. doi: 10.1126/science.1160379.
- [4] AI, Algorithmic, and Automation Incidents and Controversies, n.d. BBC presenter's AI-generated voice used to trick company. URL https://www.aiaaic.org/aiaaic-repository/ai-algorithmic-and-automation-incidents/bbc-presenters-ai-generated-voice-used-to-trick-company. Accessed: 2024-07-09.
- [5] AI21 Labs, n.d. AI21 labs concludes largest Turing Test experiment to date. URL https://www.ai21.com/blog/human-or-not-results. Accessed: 2024-06-01.
- [6] D. Allen and E. G. Weyl, 2024. The real dangers of generative AI. Journal of Democracy, 35(1):147–162.
- [7] American Association of Motor Vehicle Administrators, n.d. Driver license data verification (DLDV) service. URL https://www.aamva.org/getmedia/cb603635-3454-4331-b2c6-288d894f7fc4/AAMVA-DLDV-Overview-for-Customers.pdf. Accessed: 2024-08-04.
- [8] American Civil Liberties Union, 2023. The fight to stop face recognition technology. URL https://www.aclu.org/news/topic/stopping-face-recognition-surveillance.
- [9] B. Amin Azad, O. Starov, P. Laperdrix, and N. Nikiforakis, 2020. Web runner 2049: Evaluating third-party anti-bot services. In *Proceedings of the 17th International Conference on the Detection of Intrusions and Malware, and Vulnerability Assessment. DIMVA '20*, Lecture Notes in Computer Science, pages 135–159. Springer, 2020.
- [10] M. Anderljung, J. Barnhart, A. Korinek, J. Leung, C. O'Keefe, J. Whittlestone, S. Avin, M. Brundage, J. Bullock, D. Cass-Beggs, B. Chang, T. Collins, T. Fist, G. Hadfield, A. Hayes, L. Ho, S. Hooker, E. Horvitz, N. Kolt, J. Schuett, Y. Shavit, D. Siddarth, R. Trager, and K. Wolf, 2023. Frontier AI regulation: Managing emerging risks to public safety. URL http://arxiv.org/abs/2307.03718.
- [11] D. Anders, 2023. FCC issues cease-and-desist to halt targeted mortgage scam robocalls. CNET. URL https://www.cnet.com/personal-finance/mortgages/fcc-issues-cease-a nd-desist-to-halt-targeted-mortgage-scam-robocalls/.
- [12] anon aadhaar, 2024. Anon Aadhaar: A zero-knowledge protocol that allows Aadhaar ID owners to prove their identity in a privacy-preserving way. https://github.com/anon-aadhaar/anon-aadhaar.
- [13] Anthropic, 2023. Collective Constitutional AI: Aligning a language model with democratic input. URL https://www.anthropic.com/news/collective-constitutional-ai-aligning-a-language-model-with-public-input.
- [14] Anthropic, 2024. Measuring the persuasiveness of language models. URL https://www.anthropic.com/news/measuring-model-persuasiveness.
- [15] Apple Support, n.d. Add your identity cards to Wallet on iPhone. iPhone User Guide. URL https://support.apple.com/guide/iphone/add-identity-cards-iph9f1847064/ios. Accessed: 2024-06-02.

- [16] D. Assenmacher, L. Fröhling, and C. Wagner, 2024. You are a bot! Studying the development of bot accusations on Twitter. *Proceedings of the International AAAI Conference on Web and Social Media*, 18(1): 113–125. doi: 10.1609/icwsm.v18i1.31301.
- [17] M. Bai, J. G. Voelkel, J. Eichstaedt, and R. Willer, 2023. Artificial intelligence can persuade humans on political issues. OSF Preprints. URL https://osf.io/stakv.
- [18] B. Baker, I. Akkaya, P. Zhokhov, J. Huizinga, J. Tang, R. Sampedro, and J. Clune, 2022. Learning to play Minecraft with video PreTraining. OpenAI. URL https://openai.com/index/vpt/.
- [19] M. Bastos and D. Mercea, 2018. The public accountability of social platforms: Lessons from a study of bots and trolls in the Brexit campaign. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 376(2128):2018.00003.
- [20] C. Baum, O. Blazy, J.-H. Hoepman, A. Lehmann, A. Lysyanskaya, R. Mayrhofer, H. Montgomery, N. K. Nguyen, a. shelat, D. Slamanig, S. E. Thomsen, J. Camenisch, E. Lee, B. Preneel, S. Tessaro, and C. Troncoso, 2024. Cryptographers' feedback on the EU digital identity's ARF. URL https://files.dyne.org/eudi/cryptographers-feedback-june2024.pdf.
- [21] S. Bellovin, 2019. The early history of Usenet, part V: Authentication and norms. CircleID. URL https://circleid.com/posts/20191125_the_early_history_of_usenet_part_v_a uthentication_and_norms.
- [22] E. Ben-Sasson, A. Chiesa, C. Garman, M. Green, I. Miers, E. Tromer, and M. Virza, 2014. Zerocash: Decentralized anonymous payments from Bitcoin. In *Proceedings of the 2014 IEEE Symposium on Security and Privacy*, SP, pages 459–474. Institute of Electrical and Electronics Engineers, 2014. doi: 10.1109/SP.2014.36.
- [23] E. Ben-Sasson, I. Bentov, Y. Horesh, and M. Riabzev, 2018. Scalable, transparent, and post-quantum secure computational integrity. *IACR Cryptology ePrint Archive*, 2018:46.
- [24] P. Benson and V. Brannon, 2024. Section 230: A brief overview, IF12584. Congressional Research Service. URL https://crsreports.congress.gov/product/pdf/IF/IF12584.
- [25] L. Bernholz, H. Landemore, and R. Reich, editors, 2021. *Digital Technology and Democratic Theory*. University of Chicago Press.
- [26] N. Bitansky, R. Canetti, A. Chiesa, and E. Tromer, 2012. From extractable collision resistance to succinct non-interactive arguments of knowledge, and back again. In *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference*, pages 326–349, 2012.
- [27] R. Bloemen, D. Kales, P. Sippl, and R. Walch, 2024. Large-scale MPC: Scaling private iris code uniqueness checks to millions of users. IACR Cryptology ePrint Archive, Paper 2024/705. URL https://eprint.iacr.org/2024/705.
- [28] S. Bond, 2021. NYU researchers were studying disinformation on Facebook. The company cut them off. National Public Radio. URL https://www.npr.org/2021/08/04/1024791053/facebook-boots-nyu-disinformation-researchers-off-its-platform-and-critics-cry-f.
- [29] M. Borge, E. Kokoris-Kogias, P. Jovanovic, L. Gasser, N. Gailly, and B. Ford, 2017. Proof-of-Personhood: Redemocratizing permissionless cryptocurrencies. In *2017 IEEE European Symposium on Security and Privacy Workshops*, EUROS&PW, pages 23–26. IEEE Computer Society, 2017. doi: 10.1109/EuroSPW.2017.46.
- [30] d. boyd, 2012. The politics of "real names". *Communications of the ACM*, 55(8):29–31. doi: 10.1145/2240 236.2240247.
- [31] S. Bradshaw, H. Bailey, and P. Howard, 2021. *Industrialized Disinformation: 2020 Global Inventory of Organized Social Media Manipulation*. Oxford Internet Institute.

- [32] A. W. Branscomb, 1995. Anonymity, autonomy, and accountability: Challenges to the First Amendment in cyberspaces. *Yale Law Journal*, 104(7):1639–1679. doi: 10.2307/797027.
- [33] S. Bray, S. Johnson, and B. Kleinberg, 2023. Testing human ability to detect "deepfake" images of human faces. *Journal of Cybersecurity*, 9(1). doi: 10.1093/cybsec/tyad01.
- [34] British Columbia Public Service, n.d. Person credential. URL https://digital.gov.bc.ca/digital-trust/online-identity/person-credential/.
- [35] T. Brooks, B. Peebles, C. Holmes, W. DePue, Y. Guo, L. Jing, D. Schnurr, J. Taylor, T. Luhman, E. Luhman, C. Ng, R. Wang, and A. Ramesh, 2024. Video generation models as world simulators. OpenAI. URL https://openai.com/research/video-generation-models-as-world-simulators.
- [36] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, 2020. Language models are few-shot learners. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS '20, pages 1877–1901. Curran Associates Inc., 2020.
- [37] J. Buolamwini and T. Gebru, 2018. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Proceedings of the 1st Conference on Fairness, Accountability and Transparency*, Machine Learning Research, pages 77–91. PMLR, 2018.
- [38] V. Buterin, 2023. What do I think about biometric proof of personhood? Blog. URL https://vitalik.eth.limo/general/2023/07/24/biometric.html.
- [39] J. Calma, 2023. Twitter just closed the book on academic research. The Verge. URL https://www.theverge.com/2023/5/31/23739084/twitter-elon-musk-api-policy-chilling-academic-research.
- [40] J. Camenisch and A. Lysyanskaya, 2002. Dynamic accumulators and application to efficient revocation of anonymous credentials. In M. Yung, editor, *Proceedings of the 22nd Annual International Cryptology Conference: Advances in Cryptology. CRYPTO '02*, Lecture Notes in Computer Science. Springer, 2002. doi: 10.1007/3-540-45708-9 5.
- [41] J. Camenisch and A. Lysynskaya, 2001. An efficient system for non-transferable anonymous credentials with optional anonymity revocation. In B. Pfitzmann, editor, *Proceedings of the International Conference on the Theory and Application of Cryptographic Techniques: Advances in Cryptology. EUROCRYPT '01*, Lecture Notes in Computer Science, pages 93–118. Springer, 2001. doi: 10.1007/3-540-44987-6 7.
- [42] J. Camenisch and E. Van Herreweghen, 2002. Design and implementation of the idemix anonymous credential system. In *Proceedings of the 9th ACM Conference on Computer and Communications Security*, CCS '02, pages 21–30. Association for Computing Machinery, 2002. doi: 10.1145/586110.586114.
- [43] A. Chan, C. Ezell, M. Kaufmann, K. Wei, L. Hammond, H. Bradley, E. Bluemke, N. Rajkumar, D. Krueger, N. Kolt, L. Heim, and M. Anderljung, 2024. Visibility into AI agents. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '24, pages 958–973. Association for Computing Machinery, 2024. doi: 10.1145/3630106.3658948.
- [44] A. Chan, N. Kolt, P. Wills, U. Anwar, C. Schroeder de Witt, N. Rajkumar, L. Hammond, D. Krueger, L. Heim, and M. Anderljung, 2024. IDs for AI systems. URL https://arxiv.org/abs/2406.12137.
- [45] D. Chaum, 1981. Untraceable electronic mail, return addresses, and digital pseudonyms. *Communications of the ACM*, pages 84–90.
- [46] D. Chaum, 1983. Blind signatures for untraceable payments. In D. Chaum, R. L. Rivest, and A. T. Sherman, editors, *Proceedings of the International Cryptology Conference: Advances in Cryptology. CRYPTO '82*, pages 199–203. Plenum Press, 1983.

- [47] D. Chaum, 1985. Security without identification: Transaction systems to make Big Brother obsolete. *Communications of the ACM*, 28(10):1030–1044. doi: 10.1145/4372.4373.
- [48] D. Chaum, A. Fiat, and M. Naor, 1988. Untraceable electronic cash: (extended abstract). In S. Goldwasser, editor, *Proceedings of the 8th Annual International Cryptology Conference: Advances in Cryptology. CRYPTO* '88, volume 403 of *Lecture Notes in Computer Science*, pages 319–327. Springer, 1988. doi: 10.1007/0-387 -34799-2 25.
- [49] B. Chen, 2023. Everyone wants your email address. Think twice before sharing it. New York Times. URL https://www.nytimes.com/2023/01/25/technology/personaltech/email-address-digital-tracking.html.
- [50] H. Chen and K. Magramo, 2024. Finance worker pays out \$25 million after video call with deepfake "Chief Financial Officer". CNN. URL https://www.cnn.com/2024/02/04/asia/deepfake-cfo-scam-hong-kong-intl-hnk/index.html.
- [51] K.-P. Chen, T. Liang, T. Chang, Y.-C. Liu, S.-Y. Yin, and Y.-T. Yu, 2020. How serious is shill bidding in online auctions? Evidence from eBay Motors. SSRN. URL https://idv.sinica.edu.tw/kongpin/shill%20bid.pdf.
- [52] S. Chesterman, 2021. We, the Robots?: Regulating Artificial Intelligence and the Limits of the Law. Cambridge University Press.
- [53] R. Cheung, 2024. Mark Zuckerberg believes there could be more AI agents than people in the world. LinkedIn. URL https://www.linkedin.com/posts/rowancheung_mark-zuckerberg-believes-there-could-be-more-activity-7221939306012033026-Q3sI/.
- [54] M. Chinen, 2016. The co-evolution of autonomous machines and legal responsibility. *Virginia Journal of Law & Technology*, 20(2):338.
- [55] M. Chinen, 2019. Law and Autonomous Machines: The Co-evolution of Legal Responsibility and Technology. Edward Elgar.
- [56] S. Chopra and L. F. White, 2011. A Legal Theory for Autonomous Artificial Agents. University of Michigan Press.
- [57] P. Cihon, 2024. Chilling autonomy: Policy enforcement for human oversight of AI agents. In 41st International Conference on Machine Learning, Workshop on Generative AI and Law, GenLaw '24, 2024.
- [58] N. Clarke, 2023. Sci-fi magazine stops submissions after flood of AI-generated stories. National Public Radio. URL https://www.npr.org/2023/02/23/1159118948/sci-fi-magazine-stops-submissions-after-flood-of-ai-generated-stories.
- [59] S. Clarke and J. Savulescu, 2021. Rethinking our assumptions about moral status. In S. Clarke, H. Zohny, and J. Savulescu, editors, *Rethinking Moral Status*. Oxford University Press. URL http://www.ncbi.nlm.nih.gov/books/NBK572928/.
- [60] N. Clegg, 2024. Labeling AI-generated images on Facebook, Instagram and Threads. Meta. URL https: //about.fb.com/news/2024/02/labeling-ai-generated-images-on-facebook-instagram-and-threads/.
- [61] Coalition for Content Provenance and Authenticity, n.d. C2PA security considerations: Threat: Stripping data within C2PA manifests. URL https://c2pa.org/specifications/specifications/1.0/security/Security_Considerations.html#_threat_stripping_c2pa_manifests. Accessed: 2024-06-04.
- [62] H. Coffey, 2024. Am I falling for ChatGPT? The dark world of AI catfishing on dating apps. Independent. URL https://www.independent.co.uk/life-style/ai-catfishing-dating-apps-chatgpt-b2531460.html.

- [63] Collective Intelligence Project, 2023. Alignment Assemblies: 2023 roadmap. URL https://cip.org/alignmentassemblies.
- [64] Committee on Homeland Security and Government Affairs, Permanent Subcommittee on Investigations, 2019. *Abuses of the Federal Notice-and-Comment Rulemaking Process: Staff Report.* United States Senate.
- [65] E. Cooke and F. Jahanian, 2005. The zombie roundup: Understanding, detecting, and disrupting botnets. In *Proceedings of the Steps to Reducing Unwanted Traffic on the Internet Workshop*, SRUTI '05, pages 39–44. USENIX Association, 2005.
- [66] J. Cox, 2024. ID verification service for TikTok, Uber, X exposed driver licenses. 404 Media. URL https://www.404media.co/id-verification-service-for-tiktok-uber-x-exposed-driver-licenses-au10tix.
- [67] J. Cox, 2024. Inside the underground site where "neural networks" churn out fake IDs. 404 Media. URL https://www.404media.co/inside-the-underground-site-where-ai-neural-networks-churns-out-fake-ids-onlyfake/.
- [68] Creator Assertions Working Group, 2024. Identity assertion: Creator assertions working group. URL https://creator-assertions.github.io/identity/1.0-draft/. Accessed: 2024-06-04.
- [69] A. Dar, 2022. We checked Elon Musks's claims about Twitter bots; here's what we found. CPO Magazine. URL https://www.cpomagazine.com/cyber-security/we-checked-elon-musks-claims-about-twitter-bots-heres-what-we-found/.
- [70] DataDome, n.d. CAPTCHA farms & challenges of CAPTCHA bot detection. URL https://datadome.co/guides/captcha/how-to-detect-captcha-farms-and-block-captcha-bots/. Accessed: 2024-07-17.
- [71] E. David, 2024. OpenAI is adding new watermarks to DALL-E 3. The Verge. URL https://www.theverge.com/2024/2/6/24063954/ai-watermarks-dalle3-openai-content-credentials.
- [72] Decentralist.com, 2023. Proof of personhood project round-up. Coinmonks. URL https://medium.com/coinmonks/proof-of-personhood-project-round-up-85925a2a2a4.
- [73] C. DerMarkar, 2022. Facilitation and traveler identification programme: The role of the International Civil Aviation Organization. ICAO Security & Facilitation. URL https://www.icao.int/MID/Documents/2022/FAL%20Webinar/PPT%203.1%20-%20MID%20FAL%20seminar-final.pdf.
- [74] M. Diamantis, 2023-01. Vicarious liability for AI. Indiana Law Journal, 99(1):317-334.
- [75] W. Diffie and M. E. Hellmen, 1976. New directions in cryptography. *ISEE Transactions on Information Theory*, IT-22(6):644–654.
- [76] R. DiResta, K. Shaffer, B. Ruppel, D. Sullivan, and R. Matney, 2019. *The Tactics & Tropes of the Internet Research Agency*. New Knowledge.
- [77] C. Dixon, 2020. Ditching free trial will cost Disney+ later this year. nScreenMedia. URL https://nscreenmedia.com/ditching-free-trials-bad-move-disney-plus/.
- [78] P. Druschel and M. F. Kaashoek, editors, 2002. First International Workshop on Peer-to-Peer Systems, Revised Papers, IPTPS '02. Springer-Verlag. doi: 10.5555/646334. URL https://archive.org/details/peertopeersystem0000iptp/page/250/mode/2up.
- [79] C. Duffy and B. Fung, 2022. Elon Musk commissioned this bot analysis in his fight with Twitter. Now it shows what he could face if he takes over the platform. CNN. URL https://www.cnn.com/2022/10/10/tech/elon-musk-twitter-bot-analysis-cyabra/index.html.

- [80] S. Dulai, 2023. Sony completes second round of AP testing of C2PA in-camera authenticity technology. Digital Photography Review. URL https://www.dpreview.com/news/9855773515/sony-associated-press-test-in-camera-authenticity-technology.
- [81] Z. Durumeric, J. Kasten, M. Bailey, and J. A. Halderman, 2013. Analysis of the HTTPS certificate ecosystem. In *Proceedings of the 2013 Internet Measurement Conference*, IMC '13, pages 291–304. Association for Computing Machinery, 2013. doi: 10.1145/2504730.2504755.
- [82] C. Dwork, J. Lotspiech, and M. Naor, 1996. Digital signets: Self-enforcing protection of digital information (preliminary version). In *Proceedings of the 28th Annual ACM Symposium on Theory of Computing*, STOC '96, page 489–498. Association for Computing Machinery, 1996. doi: 10.1145/237814.237997.
- [83] P. Eckersley, 2010. How unique is your web browser? In M. Atallah and N. Hopper, editors, *Proceedings of the 10th International Conference on Privacy Enhancing Technologies. PETS '10*, Lecture Notes in Computer Science, pages 1–18. Springer, 2010. doi: 10.1007/978-3-642-14527-8_1.
- [84] F. Eiras, A. Petrov, B. Vidgen, C. Witt, F. Pizzati, K. Elkins, S. Mukhopadhyay, A. Bibi, B. Csaba, F. Steibel, F. Barez, G. Smith, G. Guadagni, J. Chun, J. Cabot, J. M. Imperial, J. Nolazco-Flores, L. Landay, M. Jackson, P. Röttger, P. H. Torr, T. Darrell, Y. S. Lee, and J. Foerster, 2024. Near to mid-term risks and opportunities of open-source generative AI. URL https://arxiv.org/abs/2404.17047.
- [85] S. El-Sayed, C. Akbulut, A. McCroskery, G. Keeling, Z. Kenton, Z. Jalan, N. Marchal, A. Manzini, T. Shevlane, S. Vallor, D. Susser, M. Franklin, S. Bridgers, H. Law, M. Rahtz, M. Shanahan, M. H. Tessler, A. Douillard, T. Everitt, and S. Brown, 2024. A mechanism-based approach to mitigating harms from persuasive generative AI. URL https://arxiv.org/abs/2404.15058.
- [86] T. Eloundou and T. Lee, 2024. Democratic inputs to AI grant program: Lessons learned and implementation plans. OpenAI. URL https://openai.com/index/democratic-inputs-to-ai-grant-program-update/.
- [87] A. Etzioni, 2016. Keeping AI legal. Vanderbuilt Journal of Entertainment & Technology Law, 19(1):133–146. doi: 10.2139/ssrn.2726612.
- [88] eu-digital-identity wallet, 2024. European Digital Identity Wallet Architecture and Reference Framework. URL https://github.com/eu-digital-identity-wallet/eudi-doc-architecture-and-reference-framework/blob/main/docs/arf.md.
- [89] European Parliament and Council of the European Union, 2016. General data protection regulation (GDPR). Article 5(1)(c) Data Minimization. https://eur-lex.europa.eu/eli/reg/2016/679/oj.
- [90] European Telecommunications Standards Institute (ETSI), July 2024. ETSI TR 119 476: Electronic Signatures and Trust Infrastructures (ESI); Analysis of selective disclosure and zero-knowledge proofs applied to electronic attestation of attributes. URL https://www.etsi.org/deliver/etsi_tr/119400_119499/119476/01.02.01_60/tr_119476v010201p.pdf.
- [91] Europen Union Agency for Cybersecurity, 2024. Remote ID proofing good practices. URL https://www.enisa.europa.eu/publications/remote-id-proofing-good-practices.
- [92] Executive Order on the Safe, Secure, and Trustworthy Development and Use of Artificial Intelligence, 2023. Exec. Order 14110. 88 FR 75191.
- [93] F5, 2017. How cybercriminals bypass CAPTCHA. URL https://www.f5.com/company/blog/how-cybercriminals-bypass-captcha.
- [94] M. Fang, N. Damer, F. Kirchbuchner, and A. Kuijper, 2020. Demographic bias in presentation attack detection of iris recognition systems. In *Proceedings of the 28th European Signal Processing Conference*, EUSIPCO '21, pages 835–839. Institute of Electrical and Electronics Engineers, 2020. doi: 10.23919/Eusipco47968.2020 .9287321.

- [95] R. Fang, R. Bindu, A. Gupta, and D. Kang, 2024. LLM agents can autonomously hack websites. URL https://arxiv.org/abs/2402.06664.
- [96] R. Fang, R. Bindu, A. Gupta, Z. Quisi, and D. Kang, 2024. LLM agents can autonomously exploit one-day vulnerabilities. URL https://arxiv.org/abs/2404.08144.
- [97] K. Feng, N. Ritchie, P. Blumenthal, A. Parsons, and A. Zhang, 2023. Examining the impact of provenance-enabled media on trust and accuracy perceptions. *Proceedings of the ACM on Human-Computer Interaction*, 7:1–42. doi: 10.1145/3610061.
- [98] S. Fish, P. Gölz, D. C. Parkes, A. D. Procaccia, G. Rusak, I. Shapira, and M. Wüthrich, 2023. Generative social choice. URL https://arxiv.org/abs/2309.01291v2.
- [99] B. Ford and J. Strauss, 2008. An offline foundation for online accountable pseudonyms. In *Proceedings of the 1st Workshop on Social Network Systems*, SocialNets '08, pages 31–36. Association for Computing Machinery, 2008. doi: 10.1145/1435497.1435503.
- [100] K. B. Forrest, 2024. The ethics and challenges of legal personhood for AI. Yale Law Journal, 133.
- [101] Fraley v. Facebook, Inc., 2011. 830 F. Supp. 2d 785, 799 (N.D. Cal.). URL https://casetext.com/case/fraley-v-facebook-inc.
- [102] L. Franceschi-Bicchierai, 2023. 23andMe confirms hackers stole ancestry data on 6.9 million users. TechCrunch. URL https://techcrunch.com/2023/12/04/23andme-confirms-hackers-stole-ancestry-data-on-6-9-million-users/.
- [103] B. Fung, 2021. Why did OnlyFans ban sexually explicit content? It says it's the credit card companies. CNN. URL https://www.cnn.com/2021/08/20/tech/onlyfans-explicit-content-ban-payment/index.html.
- [104] A. Gabizon, 2016. How transactions between shielded addresses work. Electric Coin Co. URL https://electriccoin.co/blog/zcash-private-transactions/.
- [105] H. Gao, J. Hu, C. Wilson, Z. Li, Y. Chen, and B. Y. Zhao, 2010. Detecting and characterizing social spam campaigns. In *Proceedings of the 17th ACM Conference on Computer and Communications Security*, CCS '10, pages 681–683. Association for Computing Machinery, 2010. doi: 10.1145/1866307.1866396.
- [106] Y. Gao, G. Xu, L. Li, X. Luo, C. Wang, and Y. Sui, 2022. Demystifying the underground ecosystem of account registration bots. In *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, ESEC/FSE 2022, pages 897–909. Association for Computing Machinery, 2022.
- [107] B. Garner, 2024. Person. Black's Law Dictionary.
- [108] C. Gartenberg, 2017. Hacker beats Galaxy S8 iris scanner using an IR image and a contact lens. The Verge. URL https://www.theverge.com/circuitbreaker/2017/5/23/15680268/hacker-galaxy-s8-iris-scanner-ir-image-contact-lens-starbug.
- [109] M. Gavigan, 2022. Measuring authentication friction early in the customer journey. Deduce. URL https://www.deduce.com/measuring-authentication-friction-early-in-the-customer-journey/.
- [110] T. Germain, 2023. Meta and X are destroying their apps with paid verification. Tech News. URL https://gizmodo.com/meta-and-x-are-destroying-their-apps-with-paid-verifica-1850953400.

- [111] Global Research & Analysis Team, SecureList, 2024. Social engineering for open-source supply chain attack profit—part 2: Assessing the Y, and how, of the XZ Utils incident (social engineering). Kaspersky. URL https://securelist.com/xz-backdoor-story-part-2-social-engineering/112476/.
- [112] J. A. Goldstein, G. Sastry, M. Musser, R. DiResta, M. Gentzel, and K. Sedova, 2023. Generative language models and automated influence operations: Emerging threats and potential mitigations. URL https://arxiv.org/pdf/2301.04246.
- [113] S. Goldwasser, S. Micali, and C. Rackoff, 1989. The knowledge complexity of interactive proof systems. *SIAM Journal on Computing*, 18(1):186–208. doi: 10.1137/0218012.
- [114] E. Goodman, 2020. Digital Information Fidelity and Friction: Crafting a Systems-level Approach to Transparency. Knight First Amendment Institute at Columbia University.
- [115] J. Goodwin, 2020. Mastercard, Visa and Discover cut ties with Pornhub following allegations of child abuse. CNN. URL https://www.cnn.com/2020/12/14/business/mastercard-visa-discover-pornhub/index.html.
- [116] Google Wallet, n.d. Store your digital ID on your phone. URL https://wallet.google/digitalid/. Accessed: 2024-06-02.
- [117] P. Grassi, J. Fenton, E. Newton, R. Perlener, A. Regenscheid, W. Burr, J. Richer, N. Lefkovitz, J. Danker, Y.-Y. Choong, K. Greene, and M. Theofanos, 2023. 5.2.3: Use of biometrics. National Institute of Standards and Technology. URL https://pages.nist.gov/800-63-3/sp800-63b.html#biometric_use:{~}:text=5.2.3%20Use%20of%20Biometrics.
- [118] F. Günther and J. Hesse, editors, 2023. *Security Standardisation Research: Proceedings of the 8th International Conference, SSR 2023.* Lecture Notes in Computer Science. Springer. doi: 10.1007/978-3-031-30731-7.
- [119] L. Hagen, S. Neely, T. E. Keller, R. Scharf, and F. E. Vasquez, 2022. Rise of the machines? Examining the influence of social bots on a political discussion network. *Social Science Computer Review*, 40(2):264–287. doi: 10.1177/0894439320908190.
- [120] P. Herd, H. Hoynes, J. Michener, and D. Moynihan, 2023. Introduction: Administrative burden as a mechanism of inequality in policy implementation. *RSF: The Russell Sage Foundation Journal of the Social Sciences*, 9(5):1–30. doi: 10.7758/RSF.2023.9.5.01.
- [121] J. Ho and S. Ermon, 2016. Generative adversarial imitation learning. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, NIPS '16, pages 4572–4580. Curran Associates Inc., 2016.
- [122] J. Hoffmann, S. Borgeaud, A. Mensch, E. Buchatskaya, T. Cai, E. Rutherford, D. de Las Casas, L. A. Hendricks, J. Welbl, A. Clark, T. Hennigan, E. Noland, K. Millican, G. van den Driessche, B. Damoc, A. Guy, S. Osindero, K. Simonyan, E. Elsen, O. Vinyals, J. W. Rae, and L. Sifre, 2024. Training compute-optimal large language models. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22. Curran Associates Inc., 2024.
- [123] E. Horvitz, 2022. On the horizon: Interactive and compositional deepfakes. In *Proceedings of the 2022 International Conference on Multimodal Interaction*, ICMI '22, pages 653–661. Association for Computing Machinery, 2022. doi: 10.1145/3536221.3558175.
- [124] HuffPost, 2012. Mountain Dew naming campaign hijacked by infamous message board 4chan. URL https://www.huffpost.com/entry/4chan-mountain-dew_n_1773076.
- [125] W. Hutiri, O. Papakyriakopoulos, and A. Xiang, 2024. Not my voice! A taxonomy of ethical and safety harms of speech generators. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '24, pages 359–376. Association for Computing Machinery, 2024. doi: 10.1145/3630106.3658911.

- [126] Hyperledger, 2024. AnonCreds. Hyperledger AnonCreds Project. URL https://hyperledger.org/projects/anoncreds. Accessed: 06-13-2024.
- [127] IBM, 2023. What is LangChain? URL https://www.ibm.com/topics/langchain.
- [128] ICAO Security and Facilitation, n.d. ePassport basics. URL https://www.icao.int/Security/FAL/PKD/Pages/ePassport-Basics.aspx. Accessed: 2024-06-04.
- [129] Idena, 2021. Idena 2021 community report. URL https://medium.com/idena/idena-202 1-community-report-c7b2b1b262dc.
- [130] International Civil Aviation Organization, n.d. ICAO traveller identification programme 18th TRIP symposium presentations. URL https://www.icao.int/Meetings/TRIP-Symposium-2 023/Pages/Presentation.aspx. Accessed: 2024-06-02.
- [131] International Organization for Standardization, 2013. ISO/IEC 29115:2013, information technology security techniques entity authentication assurance framework. URL https://www.iso.org/standard/45138.html.
- [132] International Organization for Standardization, 2021. ISO/IEC 18013-5:2021, personal identification —ISO-compliant driving licence, part 5: Mobile driving licence (mDL) application. URL https://www.iso.org/standard/69084.html.
- [133] S. Jain, Z. Hitzig, and P. Mishkin, 2024. Contextual confidence and generative AI. URL https://arxiv.org/abs/2311.01193.
- [134] Joint Task Force, 2020. Security and privacy controls for information systems and organizations. National Institute of Standards and Technology. URL https://csrc.nist.gov/pubs/sp/800/53/r5/upd1/final.
- [135] C. R. Jones and B. K. Bergen, 2024. Does GPT-4 pass the Turing Test? URL http://arxiv.org/abs/2310.20216.
- [136] M. L. Jones, 2016. Ctrl + Z: The Right to Be Forgotten. New York University Press.
- [137] J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Radford, J. Wu, and D. Amodei, 2020. Scaling laws for neural language models. URL http://arxiv.org/abs/2001.08361.
- [138] F. B. Keller, D. Schoch, S. Stier, and J. Yang, 2020. Political astroturfing on twitter: How to coordinate a disinformation campaign. *Political Communication*, 37(2):256–280.
- [139] B. Kelly and L. Sisneros, 2020. Broadband access and the digital divides. Education Commission of the United States. URL https://www.ecs.org/wp-content/uploads/Broadband_Access_and_the_Digital_Divides-1-1.pdf.
- [140] M. Kelly, 2023. Watermarks aren't the silver bullet for AI misinformation. The Verge. URL https://www.theverge.com/2023/10/31/23940626/artificial-intelligence-ai-digital-watermarks-biden-executive-order.
- [141] M. Kinniment, L. Sato, H. Du, B. Goodrich, M. Hasin, L. Chan, L. Miles, T. Lin, H. Wijk, J. Burget, A. Ho, E. Barnes, and P. Christiano, 2024. Evaluating language-model agents on realistic autonomous tasks. URL https://arxiv.org/abs/2312.11671.
- [142] J. Kirchenbauer, J. Geiping, Y. Wen, J. Katz, I. Miers, and T. Goldstein, 2023. A watermark for large language models. In A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato, and J. Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning*, ICML '23, pages 17061–17084. PMLR, 2023.

- [143] N. Kohne, 2019. California's new bot law prohibits use of undeclared bots. Akin. URL https://www.akingump.com/en/insights/blogs/ag-data-dive/california-s-new-bot-law-prohibits-use-of-undeclared-bots.
- [144] N. Kolt, 2024. Governing AI agents. URL https://papers.ssrn.com/abstract=4772956.
- [145] A. Komo, S. D. Kominers, and T. Roughgarden, 2024. Shill-proof auctions. URL https://arxiv.org/abs/2404.00475.
- [146] M. Korir, S. Parkin, and P. Dunphy, 2022. An empirical study of a decentralized identity wallet: Usability, security, and perspectives on user control. In *Proceedings of the 18th Symposium on Usable Privacy and Security*, SOUPS '22. USENIX Association, 2022.
- [147] P. Korshunov and S. Marcel, 2018. DeepFakes: A new threat to face recognition? Assessment and detection. URL https://arxiv.org/abs/1812.08685.
- [148] J. Kosseff, 2022. The United States of Anonymous. Cornell University Press.
- [149] M. Koster, 1994. Important: Spiders, robots and web wanderers. email. URL https://web.archive.org/web/20131029200350/http://inkdroid.org/tmp/www-talk/4113.html.
- [150] S. Kreps and D. L. Kriner, 2023. The potential impact of emerging technologies on democratic representation: Evidence from a field experiment. *New Media & Society*. doi: 10.1177/14614448231160526.
- [151] M. Kumar, K. Jindal, and M. Kumar, 2021. A systematic survey on CAPTCHA recognition: Types, creation and breaking techniques. *Archives of Computational Methods in Engineering*, 29(2). doi: 10.1007/s11831-0 21-09608-4.
- [152] S. Kumar, J. Cheng, J. Leskovec, and V. Subrahmanian, 2017. An army of me: Sockpuppets in online discussion communities. In *Proceedings of the 26th International Conference on World Wide Web*, WWW '17, pages 857–866. International World Wide Web Conferences Steering Committee, 2017. doi: 10.1145/3038 912.3052677.
- [153] S. Lai and B. Tanner, 2022. Examining the intersection of data privacy and civil rights. Brookings. URL https://www.brookings.edu/articles/examining-the-intersection-of-data-privacy-and-civil-rights/.
- [154] J. A. Lanz, 2024. People are using basic AI to bypass KYC—but should you? Decrypt. URL https: //decrypt.co/216188/ai-generated-fake-id-bypass-kyc-aml-banks-crypt o-onlyfakes.
- [155] O. Laurent, 2013. Study exposes social media sites that delete photographs' metadata. British Journal of Photography. URL https://www.1854.photography/2013/03/study-exposes-social-media-sites-that-delete-photographs-metadata/.
- [156] A. Lior, 2020. AI entities as AI agents: Artificial intelligence liability and the AI respondent superior analogy. *Mitchell Hamlin Law Review*, 46(5):1043–1102.
- [157] Y.-J. Liu, X. Liu, C. Zhou, S. Zhang, N. Liu, Z. Liu, S. Wang, and M. Yu, 2022. How to evoke empathetic experience and deepen impression: The role of profile pictures in social media. preprint. URL https://doi.org/10.21203/rs.3.rs-2327366/v1.
- [158] A. Lysyanskaya, R. L. Rivest, A. Sahai, and S. Wolf, 2000. Pseudonym systems. In H. Heys and C. Adams, editors, *Proceedings of the 26th International Conference on Selected Areas in Cryptography. SAC '19*, Lecture Notes in Computer Science, pages 184–199. Springer, 2000. doi: 10.1007/3-540-46513-8_14.
- [159] S. Manveena, 2018. Aadhaar: India Supreme Court upholds controversial biometric database. CNN. URL https://www.cnn.com/2018/09/26/asia/india-aadhaar-ruling-intl/index.html.

- [160] M. Marlinspike and T. Perrin, 2016. The X3DH key agreement protocol. Signal. URL https://signal.org/docs/specifications/x3dh/.
- [161] T. Marlow, S. Miller, and J. T. Roberts, 2021. Bots and online climate discourses: Twitter discourse on President Trump's announcement of US withdrawal from the Paris Agreement. *Climate Policy*, 21(6): 765–777.
- [162] S. Matz, J. Teeny, S. Vaid, H. Peters, G. Harari, and M. Cerf, 2024. The potential of generative AI for personalized persuasion at scale. *Scientific Reports*, 14(1):4692. doi: 10.1038/s41598-024-53755-0.
- [163] E. McCallister, T. Grance, and K. Scarfone, 2010. Guide to protecting the confidentiality of personally identifying information (PII). National Institute of Standards and Technology. URL https://nvlpubs.nist.gov/nistpubs/legacy/sp/nistspecialpublication800-122.pdf.
- [164] McIntyre v. Ohio Elections Commission, 1995. Supreme Court of the United States 514 U.S. 334. URL https://www.oyez.org/cases/1994/93-986.
- [165] D. McKay, 2021. BC Government's Verifiable Credential Issuer Kit: Proof of Concept Report. Digital Identification and Authentication Council of Canada.
- [166] S. McKay and C. Tenove, 2021. Disinformation as a threat to deliberative democracy. *Political Research Quarterly*, 74(3):703–717.
- [167] P. McKenzie, 2022. The optimal amount of fraud is non-zero. Bits about Money. URL https://www.bitsaboutmoney.com/archive/optimal-amount-of-fraud/.
- [168] K. Mehta, A. Vemury, J. Prisby, and J. Finke, 2023. *Accelerate Adoption of Digital Identities on Mobile Devices: Identity Management*. National Cybersecurity Center of Excellence.
- [169] A. Menezes, P. van Oorschot, and S. Vanstone, 1997. *Handbook of Applied Cryptography: Discrete Mathematics and Its Applications*. CRC Press.
- [170] T. Merk, S. Cossar, and J. Kamalova, 2024. *Ethnographic Research of Proof of Humanity DAO (Investigación Etnográfica de Proof of Humanity DAO)*. European University Institute. doi: 10.2870/107946.
- [171] Meta AI, 2023. Stable Signature: A new method for watermarking images created by open source generative AI. Meta. URL https://ai.meta.com/blog/stable-signature-watermarking-generative-ai/.
- [172] P. R. Milgrom, 1981. Good news and bad news: Representation theorems and applications. *Bell Journal of Economics*, 12(2):380–91.
- [173] P. R. Milgrom and J. Roberts, 1986. Relying on the information of interested parties. *RAND Journal of Economics*, 17(1):18–32.
- [174] V. S. Miller, 1985. Use of elliptic curves in cryptography. In H. C. Williams, editor, *Proceedings of the 5th Annual International Cryptology Conference: Advances in Cryptology. CRYPTO '85*, Lecture Notes in Computer Science, pages 417–426. Springer, 1985. doi: 10.1007/3-540-39799-X 31.
- [175] E. Mitchell, Y. Lee, A. Khazatsky, C. D. Manning, and C. Finn, 2023. DetectGPT: Zero-shot machine-generated text detection using probability curvature. In A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato, and J. Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning*, ICML '23. JMLR.org, 2023. doi: 10.48550/arXiv.2301.11305.
- [176] Model Evaluation and Threat Research, 2023. Update on ARC's recent eval efforts: More information about ARC's evaluations of GPT-4 and Claude. URL https://metr.org/blog/2023-03-18-update-on-recent-evals/.

- [177] National Cybersecurity Center of Excellence, n.d. Digital identities mobile driver's license (mDL). URL https://www.nccoe.nist.gov/projects/digital-identities-mdl. Accessed: 2024-06-04.
- [178] National Institute of Standards and Technology, 2019. NIST releases data to help measure accuracy of biometric identification. URL https://www.nist.gov/news-events/news/2019/12/nist-releases-data-help-measure-accuracy-biometric-identification.
- [179] National Institute of Standards and Technology, 2020. Reducing risks posed by synthetic content: An overview of technical approaches to digital content & transparency: NIST AI 100-4. URL https://airc.nist.gov/docs/NIST.AI.100-4.SyntheticContent.ipd.pdf.
- [180] National Institute of Standards and Technology, 2020. A.5: Identity verification. URL https://pages.nist.gov/800-63-3-Implementation-Resources/63A/verification/.
- [181] National Institute of Standards and Technology, n.d. NIST SP 800-63 Digital Identity Guidelines: Call for comments on Initial Public Draft Revision 4. URL https://pages.nist.gov/800-63-4/. Accessed: 07-11-2024.
- [182] K. Nguyen, 2024. The cost of AI reasoning over time. sémaphore. URL https://semaphore.substack.com/p/the-cost-of-reasoning-in-raw-intelligence.
- [183] H. Nissenbaum, 2009. *Privacy in Context: Technology, Policy, and the Integrity of Social Life*. Stanford University Press. doi: 10.1515/9780804772891.
- [184] S. Njenga, 2024. Global privacy concerns prompt bans on Worldcoin's biometric data collection. Crypto News Flash. URL https://www.crypto-news-flash.com/global-privacy-concern s-prompt-bans-on-worldcoins-biometric-data-collection/.
- [185] Nomic, 2024. Announcing the release of GPT4ALL 3.0: The open-source LLM desktop app! URL https://www.nomic.ai/blog/posts/one-year-of-gpt4all.
- [186] M. Novak, 2023. 7 viral tweets about the Israel-Gaza conflict that are actually fake. Forbes. URL https://www.forbes.com/sites/mattnovak/2023/10/07/7-viral-tweets-a bout-the-israel-gaza-conflict-that-are-actually-fake/.
- [187] Office of Public Affairs, 2021. Justice Department and FTC announce first enforcement actions for violations of the Better Online Ticket Sales Act. U.S. Department of Justice. URL https://www.justice.gov/opa/pr/justice-department-and-ftc-announce-first-enforcement-actions-violations-better-online-ticket.
- [188] Office of Science and Technology Policy, 2022. Blueprint for an AI Bill of Rights. The White House. URL https://www.whitehouse.gov/ostp/ai-bill-of-rights/.
- [189] Office of Senator Brian Schatz, 2023. Schatz, Kennedy introduce bipartisan legislation to provide more transparency on AI-generated content. URL https://www.schatz.senate.gov/news/press-releases/schatz-kennedy-introduce-bipartisan-legislation-to-provide-more-transparency-on-ai-generated-content.
- [190] Office of Senator John Barrasso, 2024. Barrasso bill stops foreign and artificial intelligence influence on energy and federal lands policy. URL https://www.barrasso.senate.gov/public/index.cfm/2024/7/barrasso-bill-stops-foreign-and-artificial-intelligence-influence-on-energy-and-federal-lands-policy.
- [191] Office of the New York State Attorney General, 2021. Attorney General James issues report detailing millions of fake comments, revealing secret campaign to influence FCC's 2017 repeal of net neutrality rules. URL https://ag.ny.gov/press-release/2021/attorney-general-james-issues-report-detailing-millions-fake-comments-revealing.

- [192] P. Ohlhaver, M. Nikulin, and P. Berman, 2024. Compressed to 0: The silent strings of Proof of Personhood. *SSRN Electronic Journal*. doi: 10.2139/ssrn.4749892.
- [193] P. Olson, 2024. On Zoom, "you're on mute" is now "are you real?". Bloomberg.com. URL https://www.bloomberg.com/opinion/articles/2024-02-05/a-25-million-hong-kong-deepfake-scam-on-zoom-shows-new-ai-risks.
- [194] P. Oltermann, 2022. European leaders tricked by deepfake video calls with mayor of Kyiv. The Guardian. URL https://www.theguardian.com/world/2022/jun/25/european-leaders-deepfake-video-calls-mayor-of-kyiv-vitali-klitschko. Accessed: 2024-06-13.
- [195] OpenAI, 2023. GPT-4V(ision) system card. URL https://cdn.openai.com/papers/GPTV_S ystem_Card.pdf.
- [196] OpenAI, 2024. Navigating the challenges and opportunities of synthetic voices. URL https://openai.com/index/navigating-the-challenges-and-opportunities-of-synthetic-voices/.
- [197] OpenAI, n.d. DALL·E 2. URL https://openai.com/index/dall-e-2/. Accessed: 2024-06-04.
- [198] OpenAI et al., 2024. GPT-4 technical report. URL http://arxiv.org/abs/2303.08774.
- [199] A. Othman and J. Callahan, 2018. The horcrux protocol: A method for decentralized biometric-based self-sovereign identity. In *Proceedings of the International Joint Conference on Neural Networks*, IJCNN, pages 1–7. Institute of Electrical and Electronics Engineers, 2018. doi: 10.1109/IJCNN.2018.8489316.
- [200] PAI Staff, 2023. Building a glossary for synthetic media transparency methods, part 1: Indirect disclosure. Partnership on AI. URL https://partnershiponai.org/glossary-for-synthetic-media-transparency-methods-part-1-indirect-disclosure/.
- [201] A. Palmer and A. Spirling, 2023. Large language models can argue in convincing ways about politics, but humans dislike AI authors: Implications for governance. *Political Science*, 75(3):281—291. doi: 10.1080/00323187.2024.2335471.
- [202] M. Panditharatne, D. Weiner, and D. Kriner, 2023. Artificial intelligence, participatory democracy, and responsive government. Brenner Center for Justice. URL https://www.brennancenter.org/our-work/research-reports/artificial-intelligence-participatory-democracy-and-responsive-government.
- [203] C. Paquin and G. Zaverucha, 2013. U-Prove cryptographic specification V1.1, Revision 3. Microsoft Corporation. URL https://www.microsoft.com/en-us/research/wp-content/uploads/2016/02/U-Prove20Cryptographic20Specification20V1.1.pdf.
- [204] F. Petillion and J. Janssen, 2017. Competing for the Internet: ICAAN Gate An Analysis and Plea for Judicial Review through Arbitration. Wolters Kluwer.
- [205] M. Phuong, M. Aitchison, E. Catt, S. Cogan, A. Kaskasoli, V. Krakovna, D. Lindner, M. Rahtz, Y. Assael, S. Hodkinson, H. Howard, T. Lieberum, R. Kumar, M. A. Raad, A. Webson, L. Ho, S. Lin, S. Farquhar, M. Hutter, G. Deletang, A. Ruoss, S. El-Sayed, S. Brown, A. Dragan, R. Shah, A. Dafoe, and T. Shevlane, 2024. Evaluating frontier models for dangerous capabilities. URL http://arxiv.org/abs/2403.13793.
- [206] A. Polyakova and C. Meserole, 2019. Exporting digital authoritarianism. Brookings Institute. URL https://www.brookings.edu/wp-content/uploads/2019/08/fp_20190826_digital_authoritarianism_polyakova_meserole.pdf.
- [207] Privacy.com, n.d. What are virtual cards? URL https://privacy.com/virtual-card. Accessed: 2024-02-06.

- [208] B. Quinn, 2023. Artificial intelligence tools for detection, research and writing: AI detection. Texas Tech University Libraries. URL https://guides.library.ttu.edu/artificialintellige ncetools/detection.
- [209] A. Quito, 2018. There's a thriving black market for selfies with passports. Quartz. URL https://qz.com/1422783/the-passport-scan-scam-theres-a-thriving-black-market-for-selfies-with-pictures-of-passports.
- [210] Radware, 2016. A game of cat and mouse: Dynamic IP address and cyber attacks. URL https://radware.com/security/ddos-threats-attacks/ddos-attack-types/dynamic-ip-address-cyber-attacks/.
- [211] J. M. Rao and D. H. Reiley, 2012. The economics of spam. *Journal of Economic Perspectives*, 26(3):87–110. doi: 10.1257/jep.26.3.87.
- [212] J. Ratkiewicz, M. Conover, M. Meiss, B. Gonçalves, S. Patil, A. Flammini, and F. Menczer, 2011. Truthy: Mapping the spread of astroturf in microblog streams. In *Proceedings of the 20th International Conference Companion on World Wide Web*, WWW '11, page 249–252. Association for Computing Machinery, 2011. doi: 10.1145/1963192.1963301.
- [213] Republic of Estonia: Information System Authority, n.d. ID-card and its uses. URL https://www.id.ee/en/article/id-card-and-its-uses/. Accessed: 2024-06-06.
- [214] Reuters, 2024. Hong Kong regulator directs Worldcoin to cease operations citing privacy concerns. Reuters, May 23, 2024.
- [215] Reuters, 2024. Worldcoin to resume Kenya operations after police drop investigation. Reuters, June 20, 2024.
- [216] Reuters, 2024. Sam Altman's eye-scanning Worldcoin banned in Spain. Reuters, March 6, 2024.
- [217] R. Rivest, A. Shamir, and L. Adleman, 1978. A method for obtaining digital signatures and public-key cryptosystems. *Communications of the ACM*, 21(2):120–126. doi: 10.1145/359340.359342.
- [218] R. L. Rivest, A. Shamir, and Y. Tauman, 2001. How to leak a secret. In C. Boyd, editor, *Proceedings of the 17th Annual International Conference on the Theory and Application of Cryptology and Information Security: Advances in Cryptology. ASIACRYPT '01*, Lecture Notes in Computer Science, pages 552–565. Springer, 2001. doi: 10.1007/3-540-45682-1_32.
- [219] K. Rogers, 2016. Boaty McBoatface: What you get when you let the internet decide. New York Times. URL https://www.nytimes.com/2016/03/22/world/europe/boaty-mcboatface-what-you-get-when-you-let-the-internet-decide.html.
- [220] K. Roose, 2022. An A.I.-generated picture won an art prize. Artists aren't happy. New York Times. URL https://www.nytimes.com/2022/09/02/technology/ai-artificial-intelligence-artists.html.
- [221] K. Roose, 2023. Personalized A.I. agents are here. Is the world ready for them? New York Times. URL https://www.nytimes.com/2023/11/10/technology/personalized-ai-agent s.html.
- [222] D. Sadhya and T. Sahu, 2024. A critical survey of the security and privacy aspects of the Aadhaar framework. *Computers & Security*, 140(C):103782. doi: 10.1016/j.cose.2024.103782.
- [223] J. W. Salmond and G. Williams, 1947. Jurisprudence. Sweet & Maxwell.
- [224] V. Sankaran, 2023. Elon Musk says Twitter to move behind paywall as all users forced to pay "small" monthly fee. Independent. URL https://www.independent.co.uk/tech/elon-musk-charging-for-twitter-b2414002.html.

- [225] K. Schiff, D. Schiff, and N. Bueno, 2024. The liar's dividend: Can politicians use deepfakes and fake news to evade accountability? *American Political Science Review*, pages 1–20. doi: 10.1017/S0003055423001454.
- [226] J. Schweitzer, 2022. How to address the administrative burdens of accessing the safety net. Center for American Progress. URL https://www.americanprogress.org/article/how-to-address-the-administrative-burdens-of-accessing-the-safety-net/.
- [227] Science and Technology Directorate, 2023. News release: DHS S&T seeks solutions for privacy preserving digital credential wallets & verifiers. U.S. Department of Homeland Security. URL https://www.dhs.gov/science-and-technology/news/2023/06/22/st-seeks-solutions-privacy-preserving-digital-credential-wallets-verifiers.
- [228] B. Sebin, 2023. The BabyAGI revolution: The dawn of autonomous AI agents. LinkedIn. URL https://www.linkedin.com/pulse/babyagi-revolution-dawn-autonomous-ai-agents-burhan-sebin/.
- [229] E. Seger, N. Dreksler, R. Moulange, E. Dardaman, J. Schuett, K. Wei, C. Winter, M. Arnold, S. Ó hÉigeartaigh, A. Korinek, M. Anderljung, B. Bucknall, A. Chan, E. Stafford, L. Koessler, A. Ovadya, B. Garfinkel, E. Bluemke, M. Aird, P. Levermore, J. Hazell, and A. Gupta, 2023. Open-Sourcing Highly Capable Foundation Models: An Evaluation of Risks, Benefits, and Alternative Methods for Pursuing Open-Source Objectives. Centre for the Governance of AI.
- [230] T. Sengupta, 2024. Russian man uses AI for online dating, claims it helped him find his wife. Hindustan Times. URL https://www.hindustantimes.com/trending/russian-man-uses-a i-for-online-dating-claims-it-helped-him-find-his-wife-10170679850 6466.html.
- [231] Y. Shavit, S. Agarwal, M. Brundage, S. Adler, C. O'Keefe, R. Campbell, T. Lee, P. Mishkin, T. Eloundou, A. Hickey, K. Slama, L. Ahmad, P. McMillan, A. Beutel, A. Passos, and D. Robinson, 2023. *Practices for Governing Agentic AI Systems*. OpenAI.
- [232] A. Shostack, 2014. Threat Modeling: Designing for Security. Wiley.
- [233] A. Siarohin, S. Lathuilière, S. Tulyakov, E. Ricci, and N. Sebe, 2019. First order motion model for image animation. In *Proceedings of the 33rd Conference on Neural Information Processing Systems*, NeurIPS '19, pages 7105–7115, 2019.
- [234] D. Siddarth, S. Ivliev, S. Siri, and P. Berman, 2020. Who watches the watchmen? A review of subjective approaches for Sybil-resistance in Proof of Personhood protocols. *Frontiers in Blockchain*, 3.
- [235] Signal Messenger, 2020–2024. libsignal: Home to the Signal Protocol as well as other cryptographic primitives which make Signal possible. URL https://github.com/signalapp/libsignal.
- [236] Signal Messenger, 2024. Signal-Server: Server supporting the Signal Private Messenger applications on Android, Desktop, and iOS. URL https://github.com/signalapp/Signal-Server. AGPL-3.0. Accessed: 2024-06-17.
- [237] S. S. Silva, R. M. Silva, R. C. Pinto, and S. R. M., 2013. Botnets: A survey. *Computer Networks*, 57(2): 378–403.
- [238] L. Solum, 1992. Legal personhood for artificial intelligences. North Carolina Law Review, 70(4):415–471.
- [239] M. Sporny, D. Longley, M. Sabadello, D. Reed, O. Steele, C. Allen, and W3C, 2024. Decentralized identifiers (DIDs) v1.1. World Wide Web Consortium Editor's Draft. URL https://w3c.github.io/did-core/.
- [240] M. Sporny, T. Thibodeau Jr., I. Herman, M. B. Jones, G. Cohen, and W3C, August 2024. Verifiable credentials data model v2.0. W3C Candidate Recommendation Draft. URL https://www.w3.org/TR/vc-data-model-2.0/. Accessed: 2024-08-10.

- [241] B. Sprunt and L. Schapitl, 2024. A neurological disease stole Rep. Jennifer Wexton's voice. An AI helped her get it back. National Public Radio. URL https://www.npr.org/2024/07/25/nx-s1-5051720/jennifer-wexton-ai-speech-progressive-supranuclear-palsy.
- [242] J. Stanley, 2023. TSA shouldn't force a bad digital ID system on America. American Civil Liberties Union. URL https://www.aclu.org/news/privacy-technology/tsa-shouldnt-force-a-bad-digital-id-system-on-america.
- [243] J. Stanley and O. Akselrod, 2022. Three key problems with the government's use of a flawed facial recognition service. American Civil Liberties Union. URL https://www.aclu.org/news/privacy-techn ology/three-key-problems-with-the-governments-use-of-a-flawed-fac ial-recognition-service.
- [244] T. Stobierski, n.d. What is selfie identity verification? Persona. URL https://withpersona.com/blog/what-is-selfie-identity-verification-and-how-does-it-work. Accessed: 2024-06-04.
- [245] The Conversation and A. Susarla, 2023. A social media researcher explains the problem with Meta and Twitter's verification services. Fast Company. URL https://www.fastcompany.com/90862625/a-social-media-researcher-explains-the-problem-with-meta-and-twitters-verification-services.
- [246] The White House, 2023. FACT SHEET: Vice President Harris announces new U.S. initiatives to advance the safe and responsible use of artificial intelligence. WhiteHouse.gov. URL https://www.whitehouse.gov/briefing-room/statements-releases/2023/11/01/fact-sheet-vice-president-harris-announces-new-u-s-initiatives-to-advance-the-safe-and-responsible-use-of-artificial-intelligence/.
- [247] K. Thomas, D. McCoy, C. Grier, A. Kolcz, and V. Paxson, 2013. Trafficking fraudulent accounts: The role of the underground market in Twitter spam and abuse. In *Proceedings of the 22nd USENIX Security Symposium*, SEC '13, pages 195–210. USENIX Association, 2013.
- [248] K. Thomas, D. Iatskiv, E. Bursztein, T. Pietraszek, C. Grier, and D. McCoy, 2014. Dialing back abuse on phone verified accounts. In *Proceedings of the 2014 ACM SIGSAC Conference on Computer and Communications Security*, CCS '14, pages 465–476. Association for Computing Machinery, 2014. doi: 10.1145/2660267.26 60321.
- [249] Thorn, 2023. How hashing and matching can help prevent revictimization. URL https://www.thorn.org/blog/hashing-detect-child-sex-abuse-imagery/.
- [250] Ticketmaster, 2023. The arms race against abusive ticket scalping. URL https://blog.ticketmaster.com/ticket-scalpers/.
- [251] M. Tobin and J. Lui, 2024. China wants to start a national ID system. New York Times. URL https://www.nytimes.com/2024/07/31/business/china-national-internet-id.html?smid=url-share.
- [252] F. Torabi, G. Warnell, and P. Stone, 2018. Behavioral cloning from observation. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, IJCAI '18, pages 4950–4957. AAAI Press, 2018.
- [253] Truepic, 2022. Truepic's authenticating camera SDK recognized by TIME's best inventions 2022. Globe-Newswire News Room. URL https://www.globenewswire.com/news-release/2022/11/10/2553559/0/en/Truepic-s-Authenticating-Camera-SDK-Recognized-by-TIME-s-Best-Inventions-2022.html.
- [254] Truepic, n.d. Truepic lens authentic mobile camera. URL https://truepic.com/truepic-lens/. Accessed: 2024-08-05.

- [255] Unique Identification Authority of India, n.d.. Authentication Requesting Agency. Government of India. URL https://uidai.gov.in/en/ecosystem/authentication-ecosystem/authentication-requesting-agency.html. Accessed: 2024-06-06.
- [256] Unique Identification Authority of India, n.d.. How does the UIDAI protect the individual and their information? FAQs. Government of India. URL https://uidai.gov.in/en/contact-support/have-any-question/277-english-uk/faqs.html. Accessed: 2024-06-06.
- [257] United Nations Refugee Agency, n.d. Guidance on registration and identity management: 5.3 documentation, implementing registration within an identity management framework. URL https://www.unhcr.org/registration-guidance/chapter5/documentation/. Accessed: 2024-06-02.
- [258] U.S. Congress, 2023. AI Labeling Act of 2023, S.2691, 118th Congress, 2023. https://www.congress.gov/bill/118th-congress/senate-bill/2691/text.
- [259] U.S. Department of Homeland Security, 2021. Increasing threat of deepfake identities. U.S. Department of Homeland Security Report. URL https://www.dhs.gov/sites/default/files/public ations/increasing_threats_of_deepfake_identities_0.pdf.
- [260] U.S. Government Accountability Office, 2023. Closing the digital divide for millions of Americans without broadband. GAO Watchblog: Following the Federal Dollar. URL https://www.gao.gov/blog/closing-digital-divide-millions-americans-without-broadband.
- [261] USA.gov, 2024. How to change your name and what government agencies to notify. URL https://www.usa.gov/name-change.
- [262] J. Vijayan, 2024. Attacker social-engineered backdoor code into XZ Utils. Dark Reading. URL https://www.darkreading.com/application-security/attacker-social-engineered-backdoor-code-into-xz-utils.
- [263] L. von Ahn, M. Blum, N. J. Hopper, and J. Langford, 2003. CAPTCHA: Using hard AI problems for security. In E. Biham, editor, *Proceedings of the Annual International Conference on the Theory and Applications of Cryptographic Techniques: Advances in Cryptology. EUROCRYPT 2003*, Lecture Notes in Computer Science, pages 294–311. Springer, 2003. doi: 10.1007/3-540-39200-9_18.
- [264] L. Wang, C. Ma, X. Feng, Z. Zhang, H. Yang, J. Zhang, Z. Chen, J. Tang, X. Chen, Y. Lin, W. X. Zhao, Z. Wei, and J. Wen, 2024. A survey on large language model based autonomous agents. *Frontiers of Computer Science*, 18(186345).
- [265] Z. M. Wang, Z. Peng, H. Que, J. Liu, W. Zhou, Y. Wu, H. Guo, R. Gan, Z. Ni, J. Yang, M. Zhang, Z. Zhang, W. Ouyang, K. Xu, S. W. Huang, J. Fu, and J. Peng, 2024. RoleLLM: Benchmarking, eliciting, and enhancing role-playing abilities of large language models. URL https://arxiv.org/abs/2310.00746.
- [266] C. Ward, K. Polglase, S. Shukla, G. Mezzofiore, and T. Lister, 2020. Russian election meddling is back via Ghana and Nigeria and in your feeds. CNN. URL https://www.cnn.com/2020/03/12/world/russia-ghana-troll-farms-2020-ward/index.html.
- [267] M. Web, M.-E. Dowling, and M. Farina, 2021. *Understanding Mass Influence: A Case Study of the Internet Research Agency as a Contemporary Mass Influence Campaign*. University of Adelaide.
- [268] M. Weiss, 2020. Investigation plan: Deepfake submissions on federal public comment servers. *Public Interest Investigations*, 2020013001.
- [269] White v. Samsung Electronics America, Inc., 1992. 971 F.2d 1395, 1399 (9th Circuit). URL https://casetext.com/case/white-v-samsung-electronics-america-inc.
- [270] B. Wolford, 2018. Everything you need to know about the "right to be forgotten". GDPR.EU. URL https://gdpr.eu/right-to-be-forgotten/.

- [271] D. Woods, 2021. I was a human CAPTCHA solver. F5 Labs. URL https://www.f5.com/labs/articles/cisotociso/i-was-a-human-captcha-solver.
- [272] World Bank, 2023. Putting people at the center of digital public infrastructure (DPI): Annual report 2023 (English). World Bank Group. URL http://documents.worldbank.org/curated/en/099 647503042425828/IDU1a9d1a6be130dc148e6193181cf9d26959fb9.
- [273] Worldcoin, 2024. Just in: The Worldcoin community has grown to 6 million verified humans. X. URL https://x.com/worldcoin/status/1811028712063701451.
- [274] X, n.d. Community Notes: Challenges. URL https://communitynotes.x.com/guide/en/about/challenges. Accessed: 2024-07-29.
- [275] X Help Center, n.d. About automated account labels. URL https://help.x.com/en/using-x/automated-account-labels. Accessed: 2024-06-06.
- [276] S. Xu, G. Chen, Y.-X. Guo, J. Yang, C. Li, Z. Zang, Y. Zhang, X. Tong, and B. Guo, 2024. VASA-1: Lifelike audio-driven talking faces generated in real time. URL https://arxiv.org/abs/2404.10667.
- [277] S. Yao, J. Zhao, D. Yu, N. Du, I. Shafran, K. Narasimhan, and Y. Cao, 2023. ReAct: Synergizing reasoning and acting in language models. URL https://arxiv.org/abs/2210.03629.
- [278] S. Yelvington, 2006. Why anonymity exists and works on newspapers' web sites. Nieman Reports, 60(4).
- [279] J. York, 2011. A case for pseudonyms. Electronic Frontier Foundation. URL https://www.eff.org/deeplinks/2011/07/case-pseudonyms.
- [280] J. C. York and D. Kayyali, 2014. Facebook's "real name" policy can cause real-world harm for the LGBTQ community. Electronic Frontier Foundation. URL https://www.eff.org/es/deeplinks/2014/09/facebooks-real-name-policy-can-cause-real-world-harm-lgbtq-community.
- [281] D. Yuan, Y. Miao, N. Z. Gong, Z. Yang, Q. Li, D. Song, Q. Wang, and X. Liang, 2019. Detecting fake accounts in online social networks at the time of registrations. In *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*, CCS '19, pages 1423–1438. Association for Computing Machinery, 2019. doi: 10.1145/3319535.3363198.
- [282] W. Zaremba, A. Dhar, L. Ahmad, T. Eloundou, S. Santurkar, S. Agarwal, and J. Leung, 2023. Democratic inputs to AI. OpenAI. URL https://openai.com/index/democratic-inputs-to-ai/.
- [283] M. Zenko, 2015. Red Team: How to Succeed By Thinking Like the Enemy. Basic Books.
- [284] H. Zhang, B. L. Edelman, D. Francati, D. Venturi, G. Ateniese, and B. Barak, 2023. Watermarks in the sand: Impossibility of strong watermarking for generative models. URL http://arxiv.org/abs/2311.04378.
- [285] J. Zhang, D. Carpenter, and M. S. Ko, 2013. Online astroturfing: A theoretical perspective. In *Americas Conference on Information Systems*, 2013.
- [286] W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong, Y. Du, C. Yang, Y. Chen, Z. Chen, J. Jiang, R. Ren, Y. Li, X. Tang, Z. Liu, P. Liu, J.-Y. Nie, and J.-R. Wen, 2023. A survey of large language models. URL https://arxiv.org/abs/2303.18223.
- [287] zk-passport, 2024. Proof of passport: Generate privacy-preserving ID proof. URL https://github.com/zk-passport/proof-of-passport. Accessed: 2024-05-31.
- [288] S. Zuboff, 2020. The Age of Surveillance Capitalism: The Fight for a Human Future at the New Frontier of Power. Public Affairs.
- [289] M. Zuckerberg, 2024. Mark Zuckerberg on Llama 3.1, open source, AI agents, safety, and more. YouTube. URL https://youtu.be/Vy3OkbtUa5k?si=UlCiohlk9tNFK_8R.

A What do we mean by "trustworthy" digital interaction?

Extensive academic literature explores the concepts of trust and trustworthiness, in both online and offline settings—too many to cite here. In this paper, when we refer to "trustworthy" digital interactions, we mean interactions that deliver on parties' reasonable expectations:¹²⁴ for instance, an online ecosystem in which service providers are confident that their services are being used as intended, and in which users can participate in digital services as intended, free from fears of abuse, attack, and other harms. Deceptive activity breaches trust—it is untrustworthy.

To set and fulfill parties' reasonable expectations, various contextual details about an interaction may need to be provided or interpreted by a party to the interaction: Who is on the other side of this interaction? What might their interests be in this matter [133]? Untrustworthy or deceptive content like disinformation often depends upon giving one party a misleading impression about the context of that content. Scams, as another example, give misleading context to trick users into thinking they are engaging with a legitimate service.

Providing people with (and fulfilling) reasonable expectations of an interaction is not the only important value; it is critical to also safeguard privacy, which may be in tension with a desire from other parties to know more about the participant. One approach for balancing these conflicts is for contextual information to be calibrated to only what is necessary to carry out the interaction. 127

As deceptive actors increasingly look to AI to execute their schemes, how do we ensure that users and service providers can verify details that properly contextualize their interactions, while restricting information disclosure to what is necessary?

Few existing approaches allow for a robust yet minimal verification of a type of information that will become increasingly important as AI capabilities continue to improve and AI tools expand their reach: proof that there is a real person behind some digital activity. From this perspective on what constitutes "trustworthy" digital interaction, personhood credentials are a natural tool to pursue. We view this pursuit as part of a broader push toward anonymous credentials [47] that prove minimal claims tailored to context (e.g., a digital credential that proves "I am over 18 years old" but nothing more).

¹²⁴We do not aim to be exhaustive with this definition.

 $^{^{125}}$ Recently, a widely shared photo on X purportedly from the center of a conflict zone was in fact pulled from a video game [186].

¹²⁶ It is possible to imagine policies that lean too far in favor of one value or another: From a privacy perspective, it would be significant overreach to require strong forms of identification for any and all activity across the Internet. Likewise, a website might struggle to uphold trustworthy interactions if it is unwilling to collect or convey any information about its users to one another—though this might also be appropriate, depending on the website's purpose.

¹²⁷ This view aligns with Helen Nissenbaum's notion of privacy as *contextual integrity*, outlined in [183]. An information flow is private if it respects the contextual norms governing the interaction. More generally, this is related to the principle of "data minimization" [89].

B Relating personhood credentials to CAPTCHAs and synthetic content transparency tools

In this section, we further explore two common alternative approaches to the problem of scalable AI-powered deception online. PHCs can complement these approaches.

B.1 CAPTCHAs and other behavioral filters

The most common behavioral filter for AI systems is CAPTCHA (Completely Automated Public Turing test to tell Computers and Humans Apart), developed in the early 2000s to differentiate humans from bots by presenting challenges that were too difficult for the bots of that era [263].

There are various types of CAPTCHAs—such as recognizing obscured words or solving rotation-based puzzles—and there are also various methods to undermine them [151].

Even before AI could directly solve CAPTCHAs, malicious actors found ways to use AI systems to bypass these barriers. For example, testers discovered that GPT-4—a text-only AI model lacking "vision" capabilities—pretended to have a disability to persuade a skeptical human to complete a CAPTCHA on its behalf [198, 176]. More broadly, it is easy for automated scripts to enlist human workers to solve these tests for fractions of a cent per CAPTCHA. 128

Now, AI systems have developed capabilities that further challenge the effectiveness of CAPTCHAs in filtering out bots online. Multimodal AI, which combines vision with cognitive abilities to perform various tasks, in particular, poses a significant challenge. 129

As AI systems get closer to exhibiting a range of human-like abilities, it will get increasingly difficult to design CAPTCHAs that all people can easily solve, but AI systems cannot.

It is worth noting that for many service providers that use CAPTCHAs, the aim is not to make bot attacks impossible but to render them economically unviable. CAPTCHA deployers might consider their measures successful if they make abuse more challenging for bots, even if some bots occasionally succeed. Nevertheless, as AI becomes less expensive, the economic barriers imposed by CAPTCHAs diminish, while humans might face increased difficulty in solving more challenging tests.

Some digital services have moved beyond CAPTCHAs to employ a broader range of behavioral filters aimed at reducing bot-based deception. For instance, anomaly detection methods attempt to identify account activity patterns that appear suspiciously coordinated or unusual compared to that of the general user base. JavaScript-based browser challenges assess factors like the manner in which users access a service. ¹³¹

¹²⁸ A bot calls on these human workers when it detects it has been presented with a CAPTCHA challenge, through services like 2captcha [70]. For detailed reports, see [271] and [93].

¹²⁹ For a review of OpenAl's GPT-4V and its current abilities to pass certain forms of CAPTCHA, see [195].

¹³⁰ For one example of anomalous behavior—unusual signup patterns—from fraudulent users, see [281]. One way that such accounts get discovered is through common patterns in their activity, which may be the result of idiosyncratic details in a malicious actor's automated script. Al agents might be able to vary their practices in ways that are harder to correlate, as opposed to following a relatively static script.

¹³¹ For an overview of such challenges and commercial providers of these, see [9]. One way that bots already attempt to appear more person-like is through taking actions like moving the mouse and page-scrolling, as well as through acting on a

These advanced methods fundamentally rely on discernible differences between how people and AI systems interact with a particular website. As AI systems become harder to distinguish from people, websites may face difficult trade-offs. To significantly reduce AI-powered abuse, they might need to limit a substantial portion of people's activity—which cannot be reliably differentiated from AI-based traffic—in the process. Personhood credentials could offer an alternative that avoids this trade-off by offering a new way to identify automated activity that does not depend on outdated assumptions about AI systems' sophistication.

B.2 Synthetic content transparency tools

Synthetic content transparency tools intend to help distinguish AI-generated ("synthetic") content from content that people created without using AI. The Partnership on AI's Responsible Practices for Synthetic Media project has outlined many such tools [200], including:

- Watermarking: proactively modifying AI-generated content (visibly or invisibly) to contain interpretable signals of how content was generated or edited. 132
- Metadata (Provenance): ¹³³including information about content's origin or its edit history, potentially with cryptographic signatures for tamper-resistance (as in C2PA—the most prevalent example of this method ¹³⁴).
- <u>Fingerprinting</u>: logging AI-generated content into a database so that future content can be checked against this, generally through hashing.¹³⁵
- <u>Classifier-based detection</u>: using AI models to assess (predict) the likelihood that content was generated by AI. ¹³⁶

There are two central reasons why personhood credentials are an important complementary tool for addressing scalable AI-powered deception online.

First, while there are promising developments in synthetic content tools' adoption and efficacy, these tools are inherently limited in scope: Given that many AI models are and will likely continue to be

time delay. We expect these actions will become more natural-seeming as AI systems learn to better imitate human website behavior, through methods such as behavior cloning [252] or other advanced techniques [121], which can allow for learning complex behavior even when the behavior's objective cannot be well-specified in advance. For an example of imitating relatively complex human behavior in digital environments, see [18].

¹³² An example of a visible watermark is the color palette applied to the bottom-right of images generated by DALL-E 2 [197]. For a discussion of invisible watermarking of AI-generated images, see [171]. For further technical details on watermarking in the domain of text, see [142].

¹³³ The Partnership on AI uses the term *metadata* to refer to this mitigation, rather than the term *provenance*, which sometimes is taken to mean a wider umbrella of mitigations for telling "where content is from" [200]. *Metadata*, in contrast, refers specifically to "labeling in metadata where the content is from."

The most prominent metadata provenance group is the Coalition for Content Provenance and Authenticity (C2PA), which aims to "[address] the prevalence of misleading information online through the development of technical standards for certifying the source and history (or provenance) of media content" [61]. C2PA is supported by a number of technology and media companies, including Adobe, the BBC, Google, Microsoft, and OpenAI, and several technology companies have begun to use its standard for tracking AI-generated images [60, 71]. Similarly, makers of physical cameras, like Sony, and of software for smartphone cameras, like Truepic, are developing tools that enable C2PA (or C2PA-like) tracking so that a photo taken with their equipment would automatically be associated with a C2PA manifest [80, 253].

¹³⁵ This approach has some similarities to hashing illegal sexual material—such as CSAM—to detect it and thus reduce its spread [249].

¹³⁶ For an overview of many methods described as "classifier-based detection," see Texas Tech's list of AI detection software [208]. For an example classifier that does not require training a new model for detection, see [175].

open-weights, we expect it will not be possible to enforce the use of these tools in all models. ¹³⁷ In addition, a perennial problem with detection tools is that outputs from a model may be modified to evade detection. ¹³⁸ In the case of metadata-based provenance, one simple evasion is to detach a piece of content from its metadata. ¹³⁹ Moreover, the ultimate value of these tools depends on designs that aid end users' interpretations of the information, which are still open areas of research [97]. Despite these drawbacks, synthetic content transparency tools are valuable and worth pursuing; they just are not a sufficient solution on their own and are best pursued in concert with other mitigations.

Second, the challenge of AI-powered deception online is significantly broader than just synthetic content; PHCs focus on a complementary part of the challenge. In common domains of interest, like social media, AI-generated content is not strictly necessary for malicious actors to still carry out AI-powered deception. For instance, AI-powered accounts might choose to amplify posts from real humans that happen to support an agenda they wish to promote. Likewise, AI-powered accounts might be used to increase the clout of particular other accounts, such as by adding to their follower count or by simply liking their posts. In both cases, there is no AI-generated content for synthetic content transparency tools to detect. Understanding whether there are real people behind these accounts, on the other hand, is well-suited to mitigating AI-powered deception even in these cases. Furthermore, as AI systems become increasingly agentic—less like content-generating machines, more like full-fledged Internet users—we expect that AI systems will increasingly interact with websites in ways that are not amenable to content-based interventions, whereas PHCs will continue to be a valuable signal. (We offer further examples below of cases in which AI-powered deception necessitates going further than synthetic content transparency tools.)

Because PHCs focus on a separate portion of the challenge from synthetic content transparency tools, the methods may be strongest as complements to one another. For example, content provenance manifests could incorporate personhood credentials as one type of supported identity attestation. For creators looking to verify their involvement with a particular piece of content, this would allow them to do so without requiring that they attach their legal identity. More broadly, synthetic transparency tools and PHCs are complementary in the general cause of encouraging people to be more discerning about who and what they see online. As Internet users become more comfortable interpreting information from one set of tools and methods, it is possible that this level of literacy will extend to others.

¹³⁷ In some open-weights implementations of models, the watermarking function can merely be removed from the model's code before running [171].

¹³⁸ Even if an invisible watermark is applied to AI-generated content, adversaries may have reliable enough methods for perturbing the content so as to remove this watermark [284]. For further discussion of methods for defeating watermarks, see [140].

¹³⁹ Some websites strip metadata from images by default to protect their users' privacy [155]. Groups like C2PA are well aware of the threat that adversaries may remove metadata from a piece of content [61]. Including signed metadata can be a useful intervention despite users' ability to strip this information: For instance, in due time, users might learn to be skeptical of images that are not paired with metadata attesting to their provenance from a physical camera. For metadata with an intention of disclosing that a piece of content was AI-generated, however, the ability to strip this metadata is a larger impediment.

¹⁴⁰ In the case of botnet attacks, for instance, malicious actors enlist AI to carry out repeated abuse of a website's services—but if the bots are primarily interacting with the website through clicks and loading particular webpages, there is no content for the website to check for evidence of being AI-generated.

¹⁴¹ Building on the C2PA standards, the Creator Assertions Working Group (CAWG) is a separate organization that has built protocols that allow creators to make additional assertions about their content—most relevant to PHCs being its identity assertions [68]. Through a CAWG identity assertion, a content creator can digitally sign a C2PA manifest so that their identity is attached to the content as it moves through digital space. Including evidence of a personhood credential in such an assertion could provide strong evidence of a human's involvement in the content's creation. CAWG has begun work on supporting Verifiable Credentials—a potential backbone of PHCs—as one form of identity assertion [68].

To emphasize the ways in which synthetic content transparency is insufficient for solving some forms of AI-powered abuse, we offer a few illustrative cases in which we expect a different approach is needed.

In some cases, AI-powered abuse can occur without relying on AI-generated content. This could occur when AI leverages human-generated content (**Case 1** and **Case 2**), or when AI interacts directly with websites (**Case 3**). In other cases, AI-generated content is central to the abuse, but transparency-centric approaches do not fully resolve the harms (**Case 4**). In yet other cases, AI-generated content is beneficial, and synthetic content transparency tools may be misinterpreted and make these uses less possible, if not complemented by other solutions (**Case 5**).

Case 1: A bad actor uses AI to spread human-generated content for harassment.

Alice wants to use AI to embarrass Bob: She uses AI to periodically create new accounts on social media, which post a real photograph of Bob doing something that she knows would embarrass him. Because the photo does not originate from AI, synthetic content transparency tools do not flag it, despite AI being used to amplify its reach.

Case 2: A bad actor uses AI to curate human-generated content to push a political agenda.

Alice wants to use AI to push a particular political agenda on social media: She uses AI to automatically analyze a stream of real people's posts, with a network of accounts that repost the material AI determines best supports her agenda. Because the material in the posts does not originate from AI, synthetic content transparency tools do not flag it, despite AI being used to automate all reposting of content by Alice's accounts.

Case 3: A bad actor deploys AI to take actions on a website and search for exploits.

Alice wants to find and exploit vulnerabilities on a website to demand a ransom payment: She uses a wide range of AI agents to analyze the site's source code and to test out different sequences of unusual inputs. Because the agents interact primarily with the website through mouse-clicks and simple text-field inputs, synthetic content transparency tools are largely unhelpful in this case. Moreover, the site receives only a very small portion of content that the AI has generated; it cannot analyze the AI's private chain-of-thought in planning and executing this abuse.

Case 4: A bad actor uses AI to generate abusive content, but identifying the content as AI-generated does not resolve the abuse.

Alice wants to use AI to bully Bob by depicting him in images she knows he will find upsetting: She uses AI to create these images and then posts them to message boards she knows Bob frequents. Because the message board does not have a blanket ban on AI-generated images—and because disclosure of the images being AI-generated does not reduce the targeted harassment that Bob feels—synthetic content transparency tools do not stop this abuse, despite AI-generated content being intrinsic to it.

¹⁴² Synthetic content transparency tools will not address mouse-clicks and other forms of browser navigation, which are action-based rather than content-based. A website could attempt to look for watermarks in content submitted to its text fields, but watermarking tends to be unreliable for short pieces of text.

Case 5: AI empowers a good actor to engage with others digitally, but the good actor's AI-generated content causes them to be mistaken for a bad actor.

Bob lost the ability to speak with much fluency after a medical incident, but he was able to recover his ability to speak by using an AI tool. Now, Bob can once again speak with others over the phone—even with his own voice. Unfortunately, many companies have started to automatically flag and filter out callers they suspect are voiced by AI, given a surge of AI-powered phone scammers. Because Bob's voice is in fact AI-generated, he is frequently flagged in these settings and has a harder time when he needs to call in for customer support. Without some alternative method, Bob struggles to prove that he is a real person—using AI to help him express himself—rather than a bad actor using AI to deceive.

C Implementation choices for personhood credentials

There are a number of ways that a PHC system can achieve the requirements we outline. As a reminder, those requirements are:

1. Credential limits (1 credential per person per issuer):

- a. Issuers check one-per-person requirement at enrollment.
- b. Expiry or regular re-authentication.

2. Unlinkable pseudonymity (privacy):

- a. Minimal identifying information stored during enrollment.
- b. Minimal disclosure during usage.
- c. Unlinkability by default.

We next describe possible methods. We do not aim to give a comprehensive summary of all possible implementations here, nor do we intend to endorse any of the implementations that are discussed here. The goal of this section is simply to provide a brief overview of the range of methods that could be employed to achieve the requirements.

C.1 Methods to achieve a one-per-person credential limit

Issuers check one-per-person requirement at enrollment

There are three main methods through which an issuer can ensure that a particular user has not received a credential from them before (and thus enforce a credential limit): existing identity documents, ¹⁴³ biometric information, and Webs of Trust (social graphs). These methods can also be combined to, in some cases, achieve higher assurance.

Importantly, an issuer can choose a verification method that originates from a different source as long as the issuer trusts its reliability. For example, a nongovernmental issuer could issue PHCs built atop government IDs, rather than devising its own method for giving only one credential per person. Likewise, a government issuer could choose to issue PHCs atop a method other than an existing government ID.

¹⁴³ Such documents are sometimes referred to as "breeder documents" in the security literature.

One option for limiting the number of credentials per person from an issuer is to rely upon existing forms of ID, so long as those systems have finite limits per person. For instance, relying upon birth certificates or tax IDs can achieve finiteness (a credential limit) if these forms of ID are also finite, though they may also have edge-cases; the issuer may need a further mechanism for deduplication in those cases. Moreover, just as some forms of ID are harder to acquire (e.g., passports) than others (e.g., library cards), a PHC provides stronger claims of being a real distinct person when it is backed by an ID that is difficult to procure in large quantities.

How can an issuer authenticate the existing form of ID, so that they do not accidentally issue personhood credentials based on a spoofed document? Some forms of existing ID are already digitally authenticatable today, whereas others may need to be physically presented in person to avoid spoofing.¹⁴⁴ One challenge may be how to authenticate IDs from different jurisdictions, in a relatively uniform manner. Usefully, there are already global standards for some forms of ID (such as the United Nations' oversight of ePassport standards), with public key infrastructure built for authenticating these; third-parties can then explore layering zero-knowledge technology atop the documents.¹⁴⁵

Biometrics are another method for limiting the number of credentials per person from an issuer; these depend on measuring a part of a person (e.g., palm, iris, fingerprint) that is persistent and unique to them, and then checking against matches from other applicants. One challenge for biometric systems is affirming the integrity of hardware devices used to collect measurements. Because biometric-based systems inherently involve the processing or storing of sensitive information about people, it is important the data be handled in a transparent and privacy-preserving manner (e.g., practices like hashing and encrypting the data). Many privacy and civil liberties groups object to biometric systems, particularly without strong checks against abuse. Some researchers have proposed methods for decreasing the risks of abuse, such as by decentralizing the storage of biometric information [2, 199, 27]. One commonly cited potential benefit of biometric-based systems is in their near universality—nearly all people will have biometric indicators (like fingerprints) to enroll in a PHC system, whereas fewer may have a

¹⁴⁴ Some national IDs hold near-field communication (NFC) readable, cryptographically signed information like name, address, and date of birth as well as face images. For a description of NFC in ePassports, see [128]. This allows a holder to demonstrably confirm their ownership of the physical card in the digital domain, unlike ID cards in countries like the United States (including some REAL ID-compliant documents), which often rely upon sending a (spoofable) photo or video for digital authentication (depending on the permissions of the verifying party).

¹⁴⁵ For more on ePassports, see [128]. For an overview of how various systems come together to enable identification for cross-border movement, see [73]. For an example of layering zero-knowledge technologies atop these documents, see [287]. ¹⁴⁶ Different biometric approaches vary in their accuracy—groups such as NIST have released de-identified datasets of biometric indicators to aid with improving the accuracy of biometric methods [178].

¹⁴⁷ It can be difficult to verify that proprietary hardware designed for collecting biometric information is free from backdoors [38]. Generalized hardware, on the other hand, might be insufficiently sensitive to prevent forms of presentation attacks [108]. One approach is for a builder of proprietary hardware—such as Worldcoin, which has built Orb devices to facilitate iris recognition—to release its schematics [27]. This can increase transparency and might ultimately help other groups to manufacture such devices, but does not fully resolve the challenge of hardware integrity.

¹⁴⁸ Some biometric systems try to limit the risks of storing biometric information by not storing other sensitive data, but there are always difficulties in determining what sensitive information can be inferred through correlations in non-sensitive information. For instance, India's Aadhaar, a large-scale biometric identity system, states that they do not store sensitive information such as religion or caste alongside biometrics [256], but there are criticisms of this approach based on the possibility of linking databases [222].

¹⁴⁹ See, for instance, criticism of direct uses of biometrics by governments and via private organizations in partnership with the government [8, 243]. At the same time, regulators may be concerned about the legality of how private actors operate biometric PHC systems—Worldcoin, for instance, has been investigated in a number of jurisdictions [214, 216, 184], and some of these bans have since been overturned [215].

¹⁵⁰ Even in cases of near-universal availability of a biometric method, an issuer will need to consider potential biases in accuracy rates across demographic groups [37, 94].

certain government ID¹⁵¹ or be well-connected in a social graph.

Web-of-Trust (WoT) is yet another method for limiting the number of credentials per person: issuers depend on social graph analysis to establish that a person is real (e.g., "Are they vouched for by a person who is strongly believed to be real?") and has not already received a credential. Often, these systems are "seeded" with initial trusted users from some other method of trust (e.g., existing IDs, biometrics, or in-person "pseudonym parties") [99]. Someone without a seed-authenticated profile can become authenticated through sufficient vouches from already-trusted parties. While WoT can be a strong method for validating that a person exists, these systems can struggle to confirm uniqueness: A person may be able to get multiple credentials from multiple distinct social circles. The most fraud-resistant mechanisms for WoT may involve layering another method, such as biometrics, atop WoT, though in a way such that users still feel comfortable with privacy assurances. To aid in enforcing one-perperson credential limits, WoT issuers may also consider other signals beyond one's social graph—such as geographic consistency, transaction patterns, and device usage (all commonly used in modern anti-fraud frameworks today).

Expiry or regular re-authentication

Importantly, issuers must protect the one-credential-per-person requirement through secure re-authentication (checking whether the credential is held by the party to whom it was issued) and/or time-bounded expiry (limits on a credential's lifespan, which impact one actor's ability to maintain multiple over time). As re-authentication can be difficult to achieve in a privacy-preserving manner, ¹⁵⁴ it may be easier for some issuers to simply set tight expiration limits on credentials.

C.2 Methods to achieve unlinkable pseudonymity in practice

Minimal identifying information stored during enrollment

The amount of information that an issuer learns and stores about an enrollee will depend on the method by which the issuer achieves its per-person credential limit, and how the issuer chooses to handle lost, stolen, or otherwise compromised¹⁵⁵ credentials. Specifically, consider a scenario in which the issuer stores no information about the user after enrollment. This preserves privacy, but if the user loses or loses control of their PHC, how does the issuer determine which PHC needs to be recovered or revoked? Different systems resolve this trade-off between privacy and ease of recovery and revocability in different ways.

¹⁵¹ As noted in **Section 4**, there are around 850 million people worldwide who do not currently hold official identity documents [272].

¹⁵² Determining whether a person has already received a credential can be a more difficult challenge for WoT-based methods. One way to assess this may be to analyze "Does their social graph look similar to cases where people have been found to have improperly obtained multiple credentials?", though this will invariably be prone to errors. For a review of WoT-based approaches, see [234].

¹⁵³ For instance, one option for layering simple biometrics atop WoT may include people having a private profile picture, which is viewable only to a small subset of trusted contacts but can still be compared against other private profile pictures in a privacy-preserving way. Use of any biometrics may undermine a system's reasons for having initially selected WoT as a method, however

 $^{^{154}}$ For an overview of authentication methods and their trade-offs, see [232].

¹⁵⁵While it is difficult to prevent the transfer of PHCs, there are methods for achieving some forms of non-transferability for anonymous credentials [41, 158] and other digital assets more generally [82].

It is possible that once a PHC is issued, the issuer stores no identifying information at all about a credential holder. For instance, even if using an existing government ID as the basis for issuing a credential, the issuer can record that an ID number has been used for this purpose without recording which PHC was issued to said ID-holder. This process still allows the issuer to check for duplicated sign ups, enabling the credential limit. 157

An issuer that does not store any identifying information must carefully design their recovery and revocation process. There is a substantial body of literature in cryptography detailing the challenge of recovery for anonymous credentials, which have much in common with PHCs [40]. Beyond the partial solutions suggested in that literature, there are other methods that would violate some of the strict requirements of anonymous credentials, but may be suitable for use in PHC systems. For instance, issuers could use methods like back-up codes, security questions (which must be chosen to not be identifying), and hardware tokens. When the issuer does store identifying information for use in recovery and revocation, it should be encrypted. 159

Furthermore, issuers may decide that a less-than-foolproof recovery and revocation process is tolerable—worth the cost for the gains in privacy. In such cases, the issuer could make other design decisions that mitigate the harms from a sometimes-faulty recovery and revocation process. For instance, as discussed above, setting a shorter length of time during which a credential is valid can mitigate the harms when a holder has lost their credential. They need only wait a short period before they would need to renew their credential anyway. Issuers can also encourage credential holders to follow best practices to safeguard their credential and avoid loss or theft in the first place. Furthermore, the presence of other PHC issuers in the broader ecosystem may mitigate the harms from losing a credential—the holder with a lost credential can acquire a new one from another issuer.

Minimal disclosure during usage and unlinkable pseudonymity

By drawing on public key cryptography, ¹⁶⁰ zero-knowledge proofs, and mechanisms like cryptographic nullifiers, PHC systems can achieve minimal disclosure and unlinkable pseudonymity. There are a wide range of protocols studied in theory and applied in practice that use these cryptographic building blocks to satisfy requirements analogous to the ones we have defined. ¹⁶¹ The following is intended only as a brief sketch of how these tools could be employed in PHC systems.

Public key cryptography is likely to be a foundational building block for a PHC system, because it allows an issuer to keep track of valid credentials in a privacy-preserving way. For instance, a PHC issuer could maintain a list of public keys (each related to a valid credential), each of which has a paired secret private key. When a new person successfully enrolls in the PHC system, the issuer lets this person add exactly one public key to the list of valid keys—the private key is known only to the user enrolling.

How does a user prove they are the holder of a valid key from the issuer, without revealing which key

¹⁵⁶ Note that PHCs are thus far more private than "real name" policies [280].

 $^{^{157}}$ One must also be careful to avoid other forms of correlation—for instance, if an issuer can link a specific person to their PHC through a field like time-of-issuance.

¹⁵⁸ For an overview of authentication methods and their tradeoffs, see [232].

¹⁵⁹ For a discussion of how to protect sensitive data, see [163].

¹⁶⁰ For foundational developments in public-key cryptography, see [75, 217, 174]. For a textbook treatment, see [169].

¹⁶¹ There are several decades worth of research into the affordances of various systems for issuing anonymous credentials. For two prominent implementations of anonymous credentials, Idemix [42] and U-Prove [203] offer practical implementations for privacy-preserving authentication.

they hold?¹⁶² One way to do this is with zero-knowledge proofs.¹⁶³ As discussed at several points in the paper, zero-knowledge proofs are cryptographic protocols that enable a "prover" to convince a "verifier" of a statement's truth, without revealing any additional information beyond the validity of the statement. In the case of PHC usage, the user who holds a PHC is the prover and a service provider is the verifier. The user proves to the service provider that the statement "I hold a valid PHC" is true, without revealing which PHC—for instance, by proving "I hold a secret private key that pairs with some public key on the issuer's list."

Beyond verifying that the user possesses a valid PHC, service providers may need a method that allows them to check whether a PHC has been used with their service before—they need to create a service-specific "pseudonym" for the credential holder. (Indeed, this ability is what enables service providers to rate-limit the use of their service with a PHC, a foundational motivation for building PHC systems.) One approach to creating service-specific pseudonyms employs cryptographic nullifiers—paralleling solutions in e-cash systems that were designed to prevent double-spending, where unique identifiers ensure a digital coin is not spent more than once [46, 48]. In the context of PHCs, when a user interacts with a service, they compute a unique nullifier—a number or string—based on their credential and the service's identity. The service can store these nullifiers to track whether a credential has been used before, without learning any information about the user's credential, as the service cannot "undo" the computation to derive the specific credential. Crucially, while the service can determine if a credential has been used in a specific context, these nullifiers are context-specific and do not allow linking the user's activities across different services or contexts.

C.3 Issuers' incentives and governance

As noted throughout the paper, there are many possible issuers of personhood credentials, including trusted institutions like governments, nonprofits, consortia, and private companies. It is important to consider the incentives of a potential PHC issuer: A democratic nation issuing PHCs has different inherent incentives than an authoritarian one, and both have different inherent incentives than a nonprofit.

Ideally, a PHC issuer's incentives are aligned with the majority of its users', i.e., the issuer would like to ensure that the system runs with high integrity and smoothly facilitates trustworthy digital activity as intended. One approach to reinforcing these incentives might be to set up a global consortium of participating organizations and governments, following a multistakeholder governance model. An analog today might be the Internet Corporation for Assigned Names and Numbers (ICANN), which coordinates the maintenance and procedures of several databases related to the namespaces and numerical identifiers

¹⁶² Note that unlinkable pseudonymity would not be achieved if the service provider were to learn which key the user holds—if so, the key could be used to track the user across different service providers.

¹⁶³ Researchers in both academia and industry are actively working on the design of zero-knowledge proof protocols—many designs are efficient enough to be deployed in the real world. Different classes of zero-knowledge protocols are optimized for different requirements. Recently, there have been significant efforts focused on developing *generic zero-knowledge*, which can be used to prove almost any statement of interest. When bandwidth concerns dominate, zk-SNARKs [26] and zk-STARKs [23] tend to be the best generic zero-knowledge choices because of their short proof sizes. MPC-in-the-head and Vector-OLE approaches may have better performance profiles when system designers instead aim to minimize the time required to generate a proof. While flexible, generic zero-knowledge proofs tend to be less concretely efficient than bespoke zero-knowledge proofs optimized to prove specific statements. A notable example of a bespoke zero-knowledge proof that could be of use in personhood credential systems is a ring signature [218], which allow provers to sign messages while only disclosing that they are a member of some publicly known group, such as the group of valid credential holders, but not which specific member of the group.

¹⁶⁴ For example implementations of cryptographic nullifiers, see [104, 22].

of the Internet. 165

A related question is how an issuer decides which use-cases their personhood credential supports. For instance, an issuer might wish to have a registration process for service providers by which they can be approved and gain access to a verification protocol, ¹⁶⁶ though this may be in tension with the requirement that issuers not gain access to usage data for any particular person's personhood credential. ¹⁶⁷ A more interoperable PHC may have benefits of reduced friction for users, who would not need to go through issuance for many different applications—though this may also increase the risk to users if they lose control of their personhood credential, and may increase the risks to privacy and civil liberties associated with a centralized issuer, detailed in Appendix D. Notably, if there are multiple issuers in a PHC ecosystem, then users can select into credential issuers with their desired level of interoperability.

D Ecosystem design and management

Recall that a personhood credentialing ecosystem has two goals: to reduce the harms of scaled deception while protecting privacy and civil liberties. We begin with a discussion of how an ecosystem in which each person can hold a bounded number of credentials (greater than one but not too many) resolves tensions in key design goals. Then we discuss service provider and issuer choices that must be considered in an ecosystem in which users hold may more than one credential.

D.1 Tensions in design goals can be resolved with bounded credentials

Here, we highlight some inherent tensions between these ecosystem design goals through two extremes, which we disfavor: one in which people can acquire an unlimited number of credentials, and one in which people can acquire only one credential from a single issuer.

Through these extremes, we illustrate that an ecosystem with bounded credentials—dependent upon having a per-person credential limit from each issuer—may best resolve these tensions. Then, we discuss how service providers might protect the integrity of their services when each person may have more than one (but not too many) PHC.

¹⁶⁵ ICANN's core governance components are a multistakeholder policy development process rolling up to a policy-making council, additional supporting organizations, and a governmental advisory committee, as well as periodic public meetings intended to field broader popular input. For a detailed study of ICANN's accountability framework and legal cases related to ICANN's review process, see [204].

¹⁶⁶ In recent years, there have been many examples of technology platforms deciding to exclude certain types of use-cases from their protocol—for instance, in the case of the credit card networks and online pornography [115, 103]. Likewise, some leading identification systems are not universal in their support of use-cases today: For instance, India's Aadhaar has an application process for certain entities to gain access to authenticating Aadhaar numbers, given a number of security requirements entailed for authenticators [255].

¹⁶⁷ It is worth noting that there are two distinct questions here: 1) Which service providers can gain access to a PHC system's verification protocol, and for what services?, and 2) When a particular verification happens, what information is made available to the issuer? By our requirements, a PHC issuer should not learn about any particular usage via verification—but the issuer might be able to limit who can access the verification protocol, without violating that requirement.

Unlimited credentials

First, consider the extreme in which issuers have no credential limits; a person can obtain many credentials from an issuer without difficulty. ¹⁶⁸ In this case, when there are unlimited credentials per person, it is hard to use these credentials to address the harms of scalable deceptive behavior. For instance, service providers cannot enforce per-person rate limits to stop repeated abuse.

On the other hand, such an ecosystem may have fewer threats to user privacy and civil liberties: Easy acquisition of credentials reduces the power of any single issuer to control the flow of credentials or their uses. Further, such an ecosystem may require less sensitive data from users: The issuer does not need to process or store any identifying information in order to check for duplication before dispensing a credential; they only need to check whether the user in question is a person.¹⁶⁹

One credential

Next, consider the opposite extreme, in which there is only one issuer of personhood credentials in the ecosystem, and each person can only obtain one personhood credential from the issuer. This approach makes it more difficult for bad actors to obtain multiple credentials for use in scaled deception, though with significant challenges to privacy and civil liberties.

One challenge is that participants will have less choice: Even if they object to the issuer's method of ensuring that people do not already have a credential, their only choice is to concede or to be excluded from participation in the ecosystem. We are concerned about these dynamics and how they may concentrate power, particularly with issuers that many people have no pre-existing reasons to trust. ¹⁷⁰¹⁷¹ More generally, having multiple issuers in an ecosystem could help encourage improvements in quality and demonstrated trustworthiness.

Bounded credentials

Given the inherent tensions highlighted by the extremes of ecosystems with unlimited credentials and ecosystems with one credential, we advocate for ecosystems that aim to bound the number of credentials obtainable by each person. This approach balances between the twin goals of addressing scaled deception and protecting privacy and civil liberties.

When each person can only obtain a bounded number of credentials in the overall ecosystem, it is feasible to enforce per-person rate limits and account creation limits. ¹⁷²

¹⁶⁸ In practice, "not having credential limits" may be a spectrum more than a binary. For instance, we expect there are significant differences between issuers that officially allow holders to have unlimited unlinked credentials, and issuers that would prefer holders to have few but do not have very effective enforcement of this.

¹⁶⁹ Such systems might also have lower stakes in the case of a holder losing their credential: If there is no limit on the number of credentials someone can obtain, a holder can be issued with a new one with fewer complications.

¹⁷⁰ There are many reasons that people might be hesitant about trusting an issuer: for instance, if they have had little contact with the issuer previously, or if the issuer is not clearly accountable in some form to the holders.

¹⁷¹ One further challenge of a single issuer is that revocation and recovery processes will have higher stakes: With only a single issuer, people would not have the ability to merely obtain a credential from another issuer instead. This might also increase an issuer's need to collect identifying information about holders so that the issuer knows which PHC to invalidate when they issue a new PHC to this person. As long as the issuer cannot also track PHC *usage* (such as via collusion with service providers), the practical benefits of preserving a link between PHCs and some aspects of a person's identity at enrollment may still outweigh the privacy costs.

¹⁷² Though it is feasible to enforce per-person rate limits, an ecosystem with a bounded number of credentials might not be

The bounded network also reduces the number of credentials in circulation, thus decreasing the supply available for sale, transfer, and theft. 173

D.2 Choices for service providers and issuers in an ecosystem with multiple issuers

When an ecosystem has multiple PHC issuers—even if each has a credential limit and thus the ecosystem has a bounded number of credentials on the whole—there are complex trade-offs that both service providers and issuers will need to consider.

For issuers, they may face decisions of whether—and in what forms—to try to cooperate with other issuers to deduplicate credential holders, who may otherwise have more than one PHC. Given that some users may have chosen a particular issuer to avoid having any relationship with a different particular issuer, under what circumstances is this appropriate?

For a service provider, decisions will often center around which issuers' PHCs to accept and how these PHC options relate to one another. Trade-offs will often entail balancing the service's ability to reduce scaled deception by ensuring each person has only one credential, compared with the implications of having fewer choices for its users (such as users feeling pressure to use a credential that is difficult for them to obtain, or the increased concentration of power in a chosen issuer, which creates challenges as discussed in **Section 4**).

There are no perfect solutions to these challenges; having multiple issuers of PHCs necessarily will involve choices of which to accept, which may affect services' abilities to rate-limit activities that are potentially deceptive. We encourage significantly more research on methods to balance these desiderata in an ecosystem. Below, we briefly discuss a number of considerations.

One principle that service providers might consider is to aim for inclusivity—trying to increase the number of their users who have access to an accepted PHC. For example, if most of a service's users live in a state that has an associated PHC backed by its state ID, the service can look to complement this by accepting a PHC that is more common among the portion who do not live in this state.

A second principle that service providers might consider is to aim for integrity—trying to reduce the number of deceptive users who would have easy access to a large number of PHCs. This may be achievable, for instance, if there are PHC systems that have non-overlapping eligibility criteria: Two PHC issuers may each require that holders of their PHC reside within their state. Because it would be difficult for a deceptive user to successfully obtain each of these PHCs at once, the service provider should expect that fewer will succeed at evading the service's per-person intentions.

We also consider these questions from a PHC issuer's perspective; are there actions that an issuer can

able to strictly achieve a *one*-per-person limit: For instance, in an ecosystem with three well-trusted PHC issuers, a service provider might decide to allow PHCs from each of these—and thus, even if the service provider allows only one verified account per PHC, someone could obtain up to three accounts by using one PHC from each issuer.

¹⁷³ Attackers face greater challenges in obtaining additional credentials if issuers do enforce per-person limits: (1) Attackers would need to manipulate the primary issuance process, which is more complex and costly because the issuer checks if they already hold a credential. (2) Purchasing credentials from a secondary market would become more expensive and difficult because there is a limited supply, and legitimate users have a higher opportunity cost for selling their sole credential. Thus, issuers adopting per-person credential limits help to enhance security against attackers without significantly hindering legitimate users from obtaining their individual credentials.

¹⁷⁴ If a user is limited to only a few PHCs—say, three instead of one—this may not be ideal for reducing deceptive behavior, but it is significantly less challenging than if a user can obtain 100.

take that would reduce the risk of duplicative use with a service, while still allowing people to have choice among PHCs?

One approach may be possible if issuers rely on the same underlying method to determine who has already received a PHC from them. For instance, if issuers each rely upon a particular state ID card, it might be possible for these issuers to coordinate on a sufficiently private mechanism for checking whether the card has previously been used with the other, prior to issuing a PHC directly to the holder.

We expect that these challenges will not be trivial to solve and may involve difficult trade-offs. Nonetheless, we believe that the benefits of an ecosystem with multiple issuers can be quite large—particularly for privacy and civil liberties—and so encourage further research and exploration.