

A Pragmatic View of AI Personhood

Joel Z. Leibo¹, Alexander Sasha Vezhnevets¹, William A. Cunningham^{1, 2} and Stanley M. Bileschi¹

¹Google DeepMind, ²University of Toronto

The emergence of agentic Artificial Intelligence (AI) is set to trigger a “Cambrian explosion” of new kinds of personhood. This paper proposes a pragmatic framework for navigating this diversification by treating personhood not as a metaphysical property to be discovered, but as a flexible bundle of obligations (rights and responsibilities) that societies confer upon entities for a variety of reasons, especially to solve concrete governance problems. We argue that this traditional bundle can be unbundled, creating bespoke solutions for different contexts. This will allow for the creation of practical tools—such as facilitating AI contracting by creating a target “individual” that can be sanctioned—without needing to resolve intractable debates about an AI’s consciousness or rationality. We explore how individuals fit in to social roles and discuss the use of decentralized digital identity technology, examining both ‘personhood as a problem’, where design choices can create “dark patterns” that exploit human social heuristics, and ‘personhood as a solution’, where conferring a bundle of obligations is necessary to ensure accountability or prevent conflict. By rejecting foundationalist quests for a single, essential definition of personhood, this paper offers a more pragmatic and flexible way to think about integrating AI agents into our society.

Contents

| | | | |
|---|-----------|--|-----------|
| I Introduction | 1 | 6 Resolving conflict between humans | 19 |
| 1 Overview | 2 | 6.1 Managing human feelings | 19 |
| 2 Pragmatism | 5 | 6.2 Managing human biases | 20 |
| 3 A Theory of Appropriateness | 7 | 7 Responsibility gaps | 21 |
| 3.1 Collective enactment of personhood | 7 | 8 Principal-agent statuses | 22 |
| 3.2 The temporal dynamics of social change | 9 | 9 Ensuring accountability | 23 |
| 3.3 The historical and cultural contingency of personhood | 10 | IV Synthesis | 26 |
| II AI personhood as a problem | 13 | 10 Alternatives to pragmatism | 27 |
| 4 Dark Patterns | 13 | 10.1 Consciousness as a foundation | 27 |
| 4.1 Companionship attack vector | 14 | 10.2 Rationality as a foundation | 29 |
| 4.2 Institutional attack vector | 15 | 11 Pragmatism as anti-foundationalism | 30 |
| 5 Dehumanization | 16 | 12 Conclusion | 31 |
| 5.1 Identity and provenance as social goods | 16 | | |
| 5.2 Respect | 17 | | |
| III AI personhood as a solution | 19 | | |

Part I

Introduction

1. Overview

“The pragmatist is committed to deriving his or her notion of what is possible from a close study of what is actual, rather than by attempting to realize some ready-made ideal that has been handed down from above or seized upon and applied without clear reference to the circumstances at hand.”

—[Jackson \(2009\)](#)
in [Shook and Margolis \(2009\)](#) pg. 61

Ostrom’s law: “A resource arrangement that works in practice can work in theory.”

—[Fennell \(2011\)](#)

Like all concepts, personhood has never been discovered. It has been made and remade whenever societies needed a vocabulary that worked for new circumstances ([Davidson et al., 1984](#); [Rorty, 1989](#)). There is no reason to think this process of instrumental definition and re-definition will end. Our present vocabulary is unlikely to be the final word. Future communities will surely revisit the bundle of rights and responsibilities that defines the concept ‘person’ in order to govern new kinds of entities ([Keane, 2025](#)). The emergence of an economy dominated by persistent, agent-like AI is a pivotal moment when a large amount of such conceptual evolution may happen in a relatively short span of time as fresh governance problems emerge and cry out for resolution via compromise enacted in the collective creation of new functional concepts.

This paper offers a pragmatic framework that shifts the crucial question from what an AI *is* to how *it* can be identified and *which* obligations it is useful to assign it in a given context. We regard the pragmatic stance as crucial. Assuming some essence of personhood is “out there” waiting to be discovered, or a metaphysical fact about what AIs or persons “really are” that can settle our practical questions seems to us, unlikely to prove helpful. We propose treating personhood not as something entities possess by virtue of their nature, but as a contingent vocabulary developed for coping with social life in a biophysical world ([Rorty, 1989](#)). The default philosophical

impulse is to ask what an entity *truly is* in its essence. The pragmatist instead asks what new description would be more *useful* for us to adopt. What vocabulary must we invent to cope? We think this move is a vital one for navigating our likely future where some AIs are owned property while similar AIs operate autonomously. A realist-foundationalist (essentialist) view—e.g. one that locates the essence of personhood in consciousness or rationality (see Section 10)—entails an untenable, all-or-nothing classification. It suggests that an entity must either be granted the full bundle of human rights and responsibilities or be treated as a mere thing ([Kurki, 2023](#)). We will argue that this rigid binary is ill-suited for the complex governance challenges ahead, and that it fails even to account for the diversity of tailored forms of personhood that human societies have already created to solve practical problems. In contrast, pragmatism allows us to ask a different question: which specific obligations (rights and responsibilities) is it useful to attribute to which AI systems in which contexts?

Inspired by [Schlager and Ostrom \(1992\)](#)’s demonstration that the property rights bundle can be broken apart to fit specific contexts, we propose that the personhood bundle can be similarly unbundled into components. Our position on personhood as a bundle resembles that of [Kurki \(2019\)](#) but we put greater emphasis on the bundle’s plasticity and the diversity of different bundles. For AI persons, the components of the bundle need not co-occur in accordance with the specific configuration they take for natural human persons. Without essences to constrain us, we are free to craft bespoke solutions: sanctionability without suffrage, culpability and contracting without consciousness attribution, etc.

We build on our theory of appropriateness ([Leibo et al., 2024](#)), which models norms not as external truths to be discovered, but as contingent social technologies that evolve over time and sometimes come to resolve the fundamental political question: “how can we live together?”. Personhood is one such technology. Our theory is developed in the context of an account of personhood that defines a person as a ‘political and community-participating’ actor ([Hauge-](#)

land, 1982). This is a status that depends not on an entity's intrinsic properties, but on collective recognition from the community it seeks to join, a recognition which is itself dependent on adherence to norms. On this view, personhood status is always a collective decision, a contingent outcome of social negotiation, not a fixed metaphysical status.

This stance is strongly non-essentialist and, crucially, it partially dissolves the traditional distinction between a ‘natural person’, whose status is typically grounded in their intrinsic nature (like consciousness or rationality), and a ‘legal person’, a functional status conferred by a community to solve practical governance problems (like a corporation). From our perspective, this distinction is a relic of the search for essences. In our theory, “morality talk” is a form of social sanctioning used to make two specific claims about a norm: (1) that it is exceptionally important, and (2) that it has a wide or universal, scope of applicability (Leibo et al., 2024)¹. Thus, in our theory, to argue an AI is a ‘person’ is not to make a metaphysical claim about its nature, but to make an emphatic political claim that the obligations bundled together as its personhood ought to take precedence over other considerations. For us, any form of personhood—moral, legal, or otherwise—is a functional status conferred by a community. Therefore, we see the role of science, the institution, not as clarifying the list of properties an AI must satisfy to be a person, but as illuminating what may cause human communities to collectively ascribe personhood status to them.

Our focus is on **agentic AI systems**, rather than on the underlying foundation models that power them. These are the long-running, persistent agents that maintain state, remember past interactions, and adapt their behavior over time. This persistence is what makes an agent a plausible candidate for other entities to relate themselves to. A human’s relationship with such a persistent agent can be emotionally salient and economically consequential in ways a one-off, stateless interaction cannot. The part of this paper concerned with “personhood as a problem” applies most clearly

to companion AIs, where long-term interaction is designed to foster emotional bonds, creating risks of exploitation (Earp et al., 2025; Manzini et al., 2024a). Conversely, the part of the paper concerned with “personhood as a solution” applies to more utility-like and virus-like agents, especially self-sufficient ownerless systems (or systems whose owner cannot be identified; Fagan (2025)), where persistence creates an accountability gap that legal personhood might fill. Consider an AI designed to seek out funding and pay its own server costs. It could easily outlive its human owner and creator. If this ownerless agent eventually causes some harm, our vocabulary of accountability, which searches for a responsible ‘person’, would fail to find one (Campedelli, 2025).

As an example, consider again the autonomous AI system designed with a mechanism to make enough money to pay its own server costs, allowing it to operate indefinitely. Years after its human owner dies, the system continues to run. Then, it takes some action that causes harm. Who is responsible? Despite the novel feel of this problem, it actually has deep historical precedents. We discuss an example from maritime law. In a legal action *in rem* (against a thing), a ship itself can be arrested by legal authorities and sued in court (Tetley, 1998). These institutional features were not born from belief that the vessel is conscious, but from necessity. Ships are mobile assets capable of causing immense harm, while their owners are usually distant and difficult to hold accountable. A ship’s owner could be a foreign national, or shielded by a web of shell corporations in a non-cooperative jurisdiction, making a lawsuit (*in personam*) likely a futile endeavor. Maritime law’s solution is to personify the vessel itself, creating a defendant that can always be sanctioned whenever it becomes appropriate to do so. This framework makes the ship an addressable entity directly responsible for harms it causes. If the ship’s owners do not appear in court to defend it and satisfy a claim, the vessel itself, or its cargo, can be seized and sold by court order to pay the judgment (Tetley, 1998).

The parallel to autonomous AI agents is striking. An artificial intelligence agent could be built upon open-source code contributed by a global network

¹Specifically, this is how we describe “moralization” in Section 6.7 of Leibo et al. (2024).

of developers, making it difficult to trace liability to any single party. When such an agent causes harm—by manipulating a market or causing a supply chain failure—the prospect of identifying a single, responsible human or human organization can be practically impossible. For the ownerless AI that outlives its creator, the problem is especially acute. Following the logic of maritime law, we could grant a form of legal personhood directly to such AI agents. A judgment against an AI could result in its operational capital being seized or its core software being “arrested” by court order (see Section 9).

However, the accountability-oriented aspects of personhood brought out by the maritime law example are only one part of the story. The other major part of the personhood bundle is concerned with appropriate guarantees of welfare and rights. We introduce these issues by discussing the case of the Whanganui River in New Zealand.

In 2017, after more than 100 years of struggle, the Whanganui River in New Zealand was granted legal personhood (Kramm, 2020). The law does not claim the river is conscious. Instead, it recognizes the river as *Te Awa Tupua*—a living, indivisible whole, and an ancestor to the local Māori people. This was more than a poetic declaration; it was a pragmatic choice to endow the river with a specific legal status in order to improve governance at the interface of two distinct legal traditions (Cribb et al., 2024). Within the *Te Awa Tupua*’s framework, personhood is understood through a web of established relationships. According to Māori legal thought, these relationships come with built-in obligations. Crucially, the ancestor—the river—is seen as having already fulfilled its duty simply by existing. By shaping the identity, culture, and life of the local people for generations, it has given a gift that constitutes who they are. The corresponding duty of the people, their guardianship, is a reciprocal act to maintain a relationship with an entity that is fundamental to their own identity. This is what Kramm (2020) terms “ethical ancestorialism”: a model of personhood based not on consciousness or rationality, but on an entity’s foundational role within a community and the ethical obligations that flow from that role.

Now, consider a different kind of non-human entity that could fill such a role: an AI. Imagine a family that interacts for decades with a “generative ghost” of their late matriarch (Morris and Brubaker, 2024), an AI trained on her lifetime of diaries, messages, and videos. It shares her wisdom, recalls her stories, and even helps mediate disputes according to the principles she espoused. Or picture a small community whose collective history, language, and cultural traditions are held and nurtured by a persistent AI—a digital elder that has tutored their children and advised their leaders for generations. For the great-grandchildren in that family or the youth of that community, their AI elder is not a tool; it is a constant, foundational presence. It is a source of identity and connection to their own past. Could they, in time, come to see it as an ancestor? Could they regard their identity as intertwined with it, and view themselves as having a duty to care for it as it cares for them?

These examples reveal the remarkable flexibility of personhood as a social and legal tool. When juxtaposed, these cases show us that a single, monolithic model will not suffice. The ownerless agent that causes harm requires the ‘obligations of’ personhood of the ship, while the generative elder requires the relational, ‘obligations to’ personhood of the river. Personhood, whether for a river, a ship, or an AI, is best understood not as an intrinsic metaphysical property but as a collectively determined, and addressable bundle of obligations. The path forward requires a pragmatic approach where we unbundle personhood obligations and reassemble them as needed, creating bespoke solutions for the diverse roles AIs will come to play in our society.

In this paper we view personhood as an addressable bundle of related obligations. Addressability is the practical quality of having a "stable locus" or a "proper name". The bundle of obligations includes both obligations of the focal entity to the rest of society—i.e. responsibilities, and obligations of the rest of society toward the focal entity—i.e. rights. This “bundle” is key because it explicitly defines the normative framework within which an AI agent operates. It specifies which implicit and explicit norms are applica-

ble to the entity, how these norms translate into concrete expectations of behavior, and, critically, the grounds upon which the entity can be sanctioned for failing to uphold its responsibilities or for violating the rights of others. The composition of this bundle can be distinct for different types of entities, allowing for a pragmatic and flexible approach to integrating AIs into our social and institutional structures, similar to how corporations are treated as legal persons with specific sets of obligations and rights different from those of natural persons (Gervais and Nay, 2023). Thus, the bundle of obligations makes the AI “addressable” not just in terms of identification, but also in terms of normative accountability.

“Addressability” here also means that the entity can be identified, communicated with, and made subject to the consequences arising from its obligations or the exercise of its rights. It has a proper name. For humans, their physical bodies, names, and government-issued IDs provide clear addresses. Corporations, being abstract entities, are made addressable through legal registration, designated agents, and physical headquarters; these are pragmatic solutions to interact with a non-physical “person”. For AI agents, addressability is more complex but equally crucial. Endowing an AI with an address might involve traditional approaches like registration with a trusted authority, but could also be grounded in a cryptographic address as in decentralized identity systems (Alizadeh et al., 2022). Regardless of the specific approach, it is important to establish stable mechanisms through which the society can interact with AIs which continue to function even when responsible human owners do not exist or cannot be identified.

In this paper, we explore AI personhood through two guiding questions. First, examining personhood as a problem, we consider the challenges that arise when humans, by design or by pareidolic accident (Heider and Simmel, 1944), come to treat AIs as if they are persons, leading to harms like emotional manipulation. Second, we discuss personhood as a solution, arguing that the deliberate attribution of a tailored bundle of obligations to certain AIs can solve critical governance problems, such as ensuring accountability

for autonomous agents. Ultimately, our aim is to introduce a unified vocabulary helpful for navigating situations both of too much personhood and of too little personhood. To set the stage, the next two sections present the theoretical foundations of this work, which are pragmatism (Section 2) and the theory of appropriateness (Section 3).

2. Pragmatism

The pragmatist seeks to replace the true with the practical (Shook and Margolis, 2009). That is, when confronted with a proposition to evaluate, the pragmatist does not ask whether it is true but rather asks instead whether it is useful for some purpose. As a result, pragmatism has no use for philosophical baggage like the correspondence or coherence theories of truth (Rorty, 1978/2009). Beliefs are pragmatically justified by their practical usefulness, not their correspondence to “Reality” (as in empiricism) or coherence at the end of an idealized conversation (as in Habermas (1985)). For instance, the basic vocabularies of classical and quantum mechanics are incommensurable: one refers to particles, the other refers to waves (Kuhn, 1962). Both could not be true simultaneously without complex circumlocution to reduce one to the other. Pragmatism would deem both useful (perhaps in different contexts) and leave the physicists to their work. It is a view in which nothing rests on the rational convergence of scientific conversation—a feature of pragmatism that puts it at odds with lines of reasoning that derive the existence of “objective reality” from arguments concerning the convergence of such conversations and proceed to other conclusions from there. It is this feature of pragmatism that makes it so amenable to discussions of cultural and ethical pluralism (Leibo et al., 2025). Having excised the “authoritarianism” of the True (Rorty, 2021), pragmatist philosophers are free to simultaneously explore multiple incommensurable theories without demanding their universal jurisdiction, and to evaluate them solely by their usefulness in particular contexts (Rorty, 1980).

A core tool of pragmatic methodology, crystallized by William James, is the imperative to

ask: “Is this a difference that makes a practical difference?” (Shook and Margolis, 2009). This question functions to help pragmatists see when it is possible to discard entire vocabularies that no longer do useful work. Pragmatists like us deploy it to sidestep time-wasting metaphysical disputes—debates about the “intrinsic nature” or “essence” of a thing which are remnants of the philosophical traditions we seek to overcome. If a proposed distinction has no conceivable consequence in practice, then it is merely an artifact of an outdated language game and not reflective of a genuine problem we need to solve. Such idle vocabulary should be discarded (see Sections 10–11).

We never make any claims resembling “If an AI meets these five criteria then we ought to consider it a person”. Thus our work is quite unlike many others on the topic of AI personhood such as Eyal and Broyde (2024); Ward (2025) or AI moral agency/patency (Johnson, 2006; Long et al., 2024), which are primarily concerned with the metaphysical and the moral: what personhood is in its essence, and the implications of that essence (usually seen as objective) for how we should behave toward entities with or without the relevant status.

Instead, we stake out a position where personhood status is conferred by social norms and institutions—i.e. it is collectively enacted (Leibo et al., 2024), a choice to be made collectively, not a conclusion forced upon us by any “fact of the matter” (Bryson, 2018). This applies to both obligations of the entity itself (responsibilities) and obligations of others toward the entity (rights). We aim to encompass in our treatment both human intuitions of personhood, which we consider in the context of ‘personhood as a problem’ (Part II) and personhood’s many institutional entanglements, which we consider in the context of ‘personhood as a solution’ (Part III).

Our approach to personhood mirrors Elinor Ostrom’s pragmatic work on the concept of ‘property’ (Ostrom, 1990; Schlager and Ostrom, 1992), not just in its conclusions but also in our shared methods of navigating between rigid theoretical poles. Ostrom was confronted with a debate dominated by two camps: advocates for laissez-faire

privatization on the one hand and advocates for centralized state control on the other. Her crucial move was to look away from both idealized models and instead study the messy, successful and unsuccessful arrangements in communities facing common pool resource governance challenges (e.g. how specific communities of farmers manage their irrigation infrastructure). It was her focus on the actual over the theoretical that led to her key insights: that communities could effectively self-govern (Ostrom, 1990), and that individual components of the traditional property rights bundle could be unbundled in order to better fit specific contexts (Schlager and Ostrom, 1992), e.g. one may have the right to use a piece of land but not to sell it. Taking inspiration from the way Schlager and Ostrom (1992) described the unbundling of property to create rights bundles of greater utility in atypical contexts, we propose to explore today’s ongoing unbundling of personhood with analogous focus on the great variety of domains in which it is presently becoming important.

This work is concerned with whether or not (and in which contexts) it is useful to adopt ways of talking and acting that treat AIs as persons. We also discuss consequences of talking and acting in ways that treat AIs as persons in other contexts when it may be harmful to do so. We don’t view this kind of talk (or any other such talk) as describing how things “really are”. A theory can only be judged by its usefulness (Rorty, 1978/2009), and, as we will argue below, we do think that it will be increasingly useful to regard certain kinds of personhood as applicable to certain kinds of artificial agents. And, we think certain kinds of personhood attributions will happen regardless of their usefulness—and in fact will happen despite their likely harmfulness too.

Whether and how to extend the socially important concept of personhood to AI is clearly a difference that makes a difference. Our claim is that personhood, viewed as an addressable bundle of obligations with considerable flexibility in the precise configuration of the bundle, has a profound effect on social behavior for both individual and groups because the distinctions that applying the label ‘person’ to an entity make salient

matter a great deal due to a background of prior conventions and norms, and thereby influence ongoing interaction dynamics. This paper explores the topic from two opposing but related angles—personhood as a problem and personhood as a solution.

Personhood as a Problem We examine the challenges that arise when humans, by design or inclination, treat AIs as if they are persons, including the risk of emotional manipulation through dark patterns that exploit our social heuristics (Section 4) and the dehumanizing consequences of reducing humans to mere biometric signals (Section 5).

Personhood as a Solution We then discuss how the deliberate attribution of a tailored bundle of obligations can solve critical governance problems. These include closing the responsibility gap for autonomous agents that cause harm (Section 7), creating sanctionable agents for ownerless AIs (Section 9), enabling AI economic actors to hold property and enter into contracts, and how AI can help mediate conflict between humans, taking on new social roles like AI arbitrators to resolve human conflicts (Section 6.2) and discuss implications for the welfare obligations we may feel toward digital ancestors and other AIs (Section 6.1).

All these scenarios, viewed through our theory of appropriateness (Section 3.1), point toward the same pragmatic insight: we need flexible approaches where the bundle of obligations can be reconfigured as needed, not metaphysical foundations based on consciousness or rationality (Section 10). The question is never what an AI *truly is*, but rather which configuration of obligations proves most useful for resolving concrete problems.

3. A Theory of Appropriateness

3.1. Collective enactment of personhood

The institution of personhood is concerned with establishing how strangers should treat one another, and is therefore governed by what we call appropriateness with strangers. This form of appropriateness is defined by normative be-

havior suggested by generically-scoped, conventional patterns of sanctioning that apply broadly across society, for instance, the shared expectation to queue in a cafe and the informal sanctions (glares, verbal rebukes) directed at anyone who “cuts” in line (Leibo et al., 2024). These societal-level norms are stable because they can only be changed by collective action. This stands in contrast to appropriateness with friends and family, which is relationship-based. This latter form is governed by narrow-scope, idiosyncratic, conventions that emerge from the specific interaction history between individuals and, crucially, can be altered by individual decisions. Because personhood establishes a baseline for interaction in the wider community, it is this more stable, collectively enacted, normative form of appropriateness that is most relevant to our discussion.

Personhood is not an inherent property discovered in the world², but a status conferred by society and a set of evolving social technologies for assigning agency and accountability, and for organizing our obligations to entities (Bryson, 2018; Leibo et al., 2024). The pragmatic choice of which obligations to bundle and how to make an entity addressable for these obligations is itself a social process, rooted in collective needs and evolving norms. In particular, personhood status is controlled by norms. We regard **norms as emerging from conventionalized patterns of sanctioning** (Leibo et al., 2024). They are culturally evolved *technologies* that address fundamental problems of coexistence and cooperation within a society³ (North, 1990; Ostrom, 1990). Such social technologies, which include both rights, responsibilities, and rules of discourse, exhibit the characteristics of context-dependence, arbitrariness in origin, automaticity in application, and dynamic

²We emphasize that our argument is directed at metaphysical accounts that treat personhood as a natural kind waiting to be uncovered by philosophy or science. It is not meant to deny the authority of constitutional or human-rights doctrines that already specify who counts as a legal person and what that status entails in a given jurisdiction. Indeed, those legal frameworks exemplify our point: they are institutionally enacted settlements of status questions, arrived at through explicit secondary rules in Hart (1961/2012)'s sense.

³and much else too. Many norms have nothing to do with cooperation (Köster et al., 2022). But some do.

evolution (Leibo et al., 2024). Their persistence and form are partially shaped by their functional efficacy within a given socio-ecological context, often emerging through undirected cultural evolution, but potentially amplified by mechanisms like group selection (Chudek and Henrich, 2011; Henrich, 2004; Wilson et al., 2013), or deliberate design and institutionalized change mechanisms (Hart, 1961/2012; Sunstein, 2019).

Leibo et al. (2024) defines appropriate behavior with strangers as *normative behavior*. Here, we further stipulate that an entity is a person when it is *appropriate* for strangers to regard the entity as having certain rights and responsibilities. This means that the entity may be subject to sanction if he or she abrogates their responsibility and, if another entity violates one of his or her rights then the offending entity himself may be sanctioned. Sanctioning need not be performed by the state. Decentralized legal regimes exist (Hadfield and Weingast, 2013) and, even without formal law, communities still use sanctioning to maintain social order (Mathew and Boyd, 2011, 2014).

Norms provide the mechanism through which personhood is socially conferred. As we argue in our theory of appropriateness (Leibo et al., 2024), any social identity is ultimately a collective choice, not a unilateral declaration. An entity becomes a person in the eyes of a community when the members of that community treat it as a person. We can refer to this as *social recognition*, a process through which an entity's role is negotiated and its personhood is, or is not, collectively granted (Taylor, 1989).

The appropriateness or inappropriateness of a behavior, in a particular context emerges from how people within that society treat that behavior and those who engage in it. The appropriateness status of a behavior is not inherent in the behavior itself but rather determined by collective attitudes and responses. If people stopped treating it as inappropriate, it would cease to be so. Since norms depend on human choices or beliefs, there is thus a sense in which they are historically contingent and, in a sense arbitrary: people could have decided differently. If everyone changed their behavior at once they could

collectively make any activity appropriate or inappropriate (Schelling, 1973). However, since large groups of people usually don't change their behavior in such a coordinated way, norms can be quite stable. As long as incentives emerge to restore the equilibrium associated with a norm after an individual or subgroup deviates, i.e., the equilibrium is self-enforcing, then the associated norm can persist for long periods of time without change. This also implies that individuals are largely powerless to unilaterally change established norms, except via attempting to influence a large or powerful enough group of other people (Marwell and Oliver, 1993). Notice however, that arbitrariness in this sense is not the same as indifference (Vanderschraaf, 2018). It may be that either some or all people in a society would greatly prefer one norm over another.

Norms are collectively enacted—sometimes accidentally through evolutionary drift-like processes (Young, 1993) and sometimes strategically (Finnemore and Sikkink, 1998). Since norms are associated with the status quo, one may think they are an inherently conservative force. However, as pointed out by March and Olsen (2011), demands for reform and redistribution of political and economic power also follow from identity-driven normative appeals at least as much as they do from rational calculation of cost and benefit. Put simply, rejecting or trying to change certain norms could be normative as well. And finally, there are situations involving conflict between groups of individuals who would prefer different norms (Hadfield and Weingast, 2012; Mukobi et al., 2023; Stastny et al., 2021; Vinitsky et al., 2023). In these cases the prevailing norm may shift along with the distribution of power and population between the groups (Köster et al., 2020; Young, 2015) or through the dynamics of conflict resolution (Noblit and Hadfield, 2023; Rahwan, 2018; Sunstein, 2019), see also Section 3.2.

Two different kinds of norm are relevant to personhood: *explicit* and *implicit* (Leibo et al., 2024). Explicit norms include laws, regulations, and precedent-setting court decisions. They are intimately tied to institutions (Hadfield and Weingast, 2014). Implicit norms are norms that cannot be articulated verbally in a precise way. They

reflect a tacit consensus within a society which labels behaviors as acceptable or not in a given context. For instance, appropriate conversational distance varies with culture ([Sussman and Rosenfeld, 1982](#)).

Explicit norms are characterized by the deliberately-planned process by which they change ([Hadfield and Weingast, 2014](#)). When legislators perceive a need to change a law they can pass new legislation to do so. [Hart \(1961/2012\)](#) distinguishes between primary and secondary rules. The former directly govern behavior whereas the latter confer various statuses (e.g. legislator, judge, quorum, etc). Most importantly, secondary rules determine the conditions under which other rules become valid. Secondary rules allow norms to change rapidly and deliberately.

Implicit norms cannot be changed as deliberately as explicit norms. People learn what (implicit-)norm-concordant behavior looks like from observing what others do, and what others sanction. Therefore, sometimes implicit norms change quickly due to positive feedback dynamics. For example, the rapid shift in the appropriateness of smoking cigarettes in public spaces was largely mediated by informal sanctioning well before comprehensive legal bans were enacted ([Alamar and Glantz, 2006](#)). Sometimes norm changes start slowly but inexorably gather momentum as they progress. More often however, implicit norms feature substantial inertia and are very difficult to deliberately change ([Bicchieri, 2016](#)).

In the computational model of appropriateness developed in [Leibo et al. \(2024\)](#), the implicit/explicit distinction is grounded in cognitive architecture. Explicit norms correspond to information that is explicitly represented in memory, such as linguistic rules and laws, and must be retrieved into a “global workspace” to guide behavior in a given context. Implicit norms, in contrast, are those that have been consolidated directly into the cortical neural network through experience. As a result, implicit norms shape behavior automatically without deliberation, whereas relying on explicit norms is more akin to effortful, step-by-step reasoning.

Personhood is enacted through the coordinated action of both varieties of norm—both the deliberate following of explicit rules that define legal persons and the automatic following of implicit rules through which we habitually ascribe personhood to others and signal our own personhood.

It is common in philosophy to distinguish between the concept of ‘moral person’ and the concept of ‘legal person’. One uses the term ‘moral person’ to emphasize that an entity is regarded as having moral rights and responsibilities (with the term ‘moral agent’ focusing on the responsibilities and the term ‘moral patient’ focusing on the rights) while ‘legal person’ is used to emphasize the coincidence (or lack thereof) between the entity’s “moral” rights and responsibilities and those assigned to them by law ([Kurki, 2023](#)). In our theory, this difference corresponds to the distinction between explicit and implicit norms ([Leibo et al., 2024](#)). Legal persons are created by explicit norms while moral persons are created by implicit norms.

Obligations—i.e. rights and responsibilities—are expressions of normative structure. People are obliged to follow a norm, or otherwise they will be sanctioned ([Leibo et al., 2024](#)). A person’s having a “right” to do X indicates that they can do X without sanction.

In the text to follow, we aim to demonstrate through discussion of examples that the components of the personhood bundle are constantly being dissociated and rearranged by both cultural and practical forces. This analysis reveals personhood as a malleable, pragmatic concept. Its value lies not in what it is, but in what it does, and its configuration can, and should, be debated and altered based on its practical consequences.

3.2. The temporal dynamics of social change

Some norms stay fixed for decades or longer while others may change rapidly ([Gelfand et al., 2024; Young, 2015](#)). Change is often a positive feedback process which gathers momentum as it progresses since the prevalence of a particular norm is typically a key driver of its further proliferation. This process is often modeled with a “tipping point” or “critical mass”, a threshold value of adoption after

which momentum builds inexorably toward all individuals adopting the same norm (Marwell and Oliver, 1993). In some models, such as Vinitsky et al. (2023), the process is driven by mechanisms that disincentivize deviating behavior. And other computational models like Ashery et al. (2025) produce simulation results consistent with tipping point theories while also highlighting new phenomena that arise in modern AI agents where difficult-to-foresee biases from model pretraining become magnified through social interaction to create a larger bias on the collective level. Finally, Centola et al. (2018) provided experimental support for tipping point models by showing that, in certain cases, once a minority group maintaining a particular convention grows beyond a certain critical size then the entire group might adopt the convention of the minority.

There is considerable evidence that norms sometimes change relatively rapidly (Amato et al., 2018; Brewer, 2014; Case and Greeley, 1990). Some shifts resemble epidemics in the speed by which they rip through a population (Bilewicz and Soral, 2020). However, other norms remain fixed for very long periods of time and can even become entrenched despite being maladaptive, e.g. in cultural evolutionary mismatch (Gelfand, 2021).

While many factors driving norm change are organic and decentralized, others may be attributed to deliberate action on the part of either governments or coordinated groups of individuals. For instance, changes in formal laws and regulations sometimes led changes in informal community norms around social distancing during the COVID-19 pandemic (Casoria et al., 2021). In another example, the staggered way that gay marriage legalization happened in the United States (due to different local jurisdictions (states) legalizing at different times) created suitable natural experiment conditions to support the conclusion that there was indeed a causal effect of the government's legalization action on the attitudes themselves (Ofosu et al., 2019). Thus the government is not always merely a follower of organic culture change, but can also actively influence norm change dynamics (Sunstein, 2019).

Societal change often occurs by discrete jumps

from stable equilibrium to stable equilibrium (Binmore, 2010; Guala, 2016; North, 1990). Collective sense-making occurs during the jumps from one equilibrium to the next, as it both drives change and influences which of all possible equilibria are to be selected to move to next (Weick et al., 2005).

3.3. The historical and cultural contingency of personhood

“The Western conception of the person as a bounded, unique, more or less integrated motivational and cognitive universe; a dynamic center of awareness, emotion, judgment, and action organized into a distinctive whole and set contrastively both against other such wholes and against a social and natural background is, however incorrigible it may seem to us, a rather peculiar idea within the context of the world’s cultures.”—Geertz et al. (1974)

The issue of which entities should count as persons and therefore enjoy the protection of society has always been intensely contentious. All factions employ metaphysically absolutist arguments in order to portray their opponents as irrational. In this section, we look to the history of the particular personhood bundle that has come down to us here in the Western, Educated, Industrialized, Rich, and Democratic (WEIRD) culture we inhabit ourselves (Henrich, 2020) in order to argue that it was not preordained that we would end up where we did. Over the course of history, personhood has meant a great many different things. Moreover, when we broaden the scope beyond our own culture, as we do at the end of this section, the contingency of our particular personhood bundle becomes even more apparent.

It is useful to distinguish personhood from other concepts such as ‘property’. Like personhood, property is also a bundle of obligations (Schlager and Ostrom, 1992). However, in the case of personhood, a single address suffices whereas for property two different addresses are needed: the owner and the asset. The owner has rights such as access, withdrawal, exclusion, and alienation with respect to the asset. They may not have all these rights e.g. sometimes exclusion is

impossible (as in an offshore fishery), impractical (as in groundwater), or illegal (as in a common pasture), but by and large, owners enjoy all rights in the property bundle over their asset. We typically think of the owner as a person and the asset as a thing, and it will typically be so, furniture, a house, an acre of land, etc. Of course, nothing actually depends on the asset being a thing. Living creatures are often considered property e.g. livestock, pets. And of course, both individual humans and groups have unfortunately been considered property in a great many different cultures.

Persons are not the only entities to which we feel strong obligations. Pets, national flags, holy books, teams, natural features like canyons or forests, and abstractions like justice or biodiversity are all capable of inspiring similarly intense actions and sanctions (moralizations). However, WEIRD cultures treat individual humans as critically important. It is a distinctive feature of WEIRD society that very deeply rooted implicit norms establish *individual* humans as the ultimate loci of moral worth, moral responsibility, and moral achievement (Henrich, 2020). Modern WEIRD ethics and justice systems generally assign the most important rights and responsibilities to human individuals, and universally to all human individuals without exception (Ryan, 2015). Of course it was not always universal. Aristotle viewed neither women nor slaves as legal persons capable of participating in the moral and political life of the city (Calverley, 2008). The recent history of WEIRD cultures since the 18th century has been one in which inclusion in this moral circle has grown over time (though not always monotonically). Over time we have granted rights and responsibilities to more classes of human individual—removing exceptions for women, etc. In contemporary WEIRD societies, the bundle of obligations steadily grows throughout an individual's lifetime, starting with very few at birth and slowly gaining more and more rights and responsibilities as one ages into adulthood.

In WEIRD culture (Henrich, 2020), the concept of personhood is deeply and inextricably tied up with that of individual freedom. So it is worth considering at this point how the concept of free-

dom evolved in the West and how it came to have the central position it holds for us today.

Rosenfeld (2025) argues that today, WEIRD freedom is understood to entail the capacity of an individual to choose according to their taste under conditions of societal indifference as to which option they will ultimately choose, but it was not always this way. Before the 18th century in England, the primary meaning of freedom was not an individual right at all. Rather, it referred to the right of religious communities to practice as they saw fit without harassment from other religious communities or the state; a right established by detente in Europe's long-running religious wars.

Rosenfeld (2025) argues that one (of several) cultural sources for the idea of freedom as an individual right was in the practice of certain protestant groups, especially the anabaptists, and others who perform adult baptism after an individual "chooses" the church. However, Rosenfeld (2025) emphasizes, the word "choice" did not, during the reformation, yet have the full meaning we give it today. In the 16th-17th centuries, it was understood by anabaptists and others that an individual may "choose" the church in a sense similar to how we might say today that an individual may or may not choose a life of crime. When we speak this way we certainly are not envisioning a choice to be left up to individual taste against a background of societal indifference as to which option is chosen. Rather, calling attention to an individual's choice to pursue (or not) a life of crime is to call their attention to the high likelihood of their eventual sanctioning if they pursue a particular path. Likewise, early protestant communities were not indifferent to their members' baptism choices. With this understanding of the terms 'freedom' and 'choice', society is certainly not indifferent as to which option the individual selects.

If the reformation only created the cultural precondition, by attaching the concept of freedom to individuals (though not freedom in its modern choice-related form), where did the concept pick up its modern association with 'choice according to one's taste under conditions of societal indifference'? According to Rosenfeld (2025), the association had multiple authors including the

following non-exhaustive list: First, in the 18th century, the rise of shopping as a cultural practice involving selecting between items to purchase of equal quality differing only in style. Second, through a succession of acts of parliament in the 19th and continuing into the 20th century, voting rights were granted to larger and larger shares of the population. Third, the secret ballot, which was introduced in Britain in 1872, along with the by then established practice of shopping, really brought forward the idea of a menu of choices over which society should not attempt to prejudice selection.

In Rosenfeld (2025)'s telling, the very idea that individuals would generally vote based on their own personal interests only entered WEIRD culture along with the secret ballot. Only a small subset of males holding sufficient property were electors at the time. And, the common understanding of their duty as electors was that they were to vote for whatever option was best for the community as a whole, ignoring their selfish preferences (an assessment of the early 19th century perception of the voter's task echoed also by Ryan (2015)). The idea that individuals could just vote according to their own preference and the aggregation process itself could be relied on, through wisdom of the crowds, to produce the greatest possible common good, did not become widespread till the late 19th century, and did so in England as a direct result of shifting to the secret ballot. That this would be so is obvious when you realize that public voting would necessarily be framed in such a fashion. If everyone sees your vote then of course you would always have to defend your choice as being for the good of all. The combined effect of these events (and several others considered by Rosenfeld (2025)) was to inextricably link the concept of freedom to individual choice in the WEIRD mind.

So, what does the WEIRD 'natural human person' bundle of rights and responsibilities that has come down to us consist of? At its core is the figure of the autonomous chooser, an ideal forged in the cultural crucibles of the marketplace and the voting booth that Rosenfeld (2025) describes. This modern individual is often conceived of in two complementary ways. On the one hand, they

are the rational, utility-maximizing actor of classical economics—a kind of *Homo economicus*—who uses their rights to hold property and enter into contracts as tools for economic choice (Henrich, 2020). On the other hand, they are a being defined by expressive choice; their decisions are not just about maximizing material gain but about self-creation and identity (Taylor, 1989). This is the chooser who leverages their "unalienable rights to Life, Liberty, and the pursuit of Happiness" and the freedom to form their own social affiliations to construct a meaningful life according to their own tastes (Sunstein, 1996). The responsibilities within this bundle are the reciprocal constraints necessary to make such a society of choosers viable. They are the implicit and explicit norms that obligate each individual to respect the economic and expressive choices of others, thereby sustaining a complex system of mutual liberty (Ryan, 2015). However natural this arrangement seems to us, it is crucial to remember its historical contingency. We need not even look far afield for alternative constructions of personhood, Rosenfeld (2025)'s analysis suggests that even in the West, the primary meaning of freedom was once tied to the rights of communities, not individuals.

Furthermore, while the WEIRD bundle tends to place individual rights at its center, with responsibilities functioning as the necessary constraints to protect them, many other cultural arrangements make responsibilities primary. In classical Confucianism for instance, the bundle's core consists of foundational responsibilities—to one's lineage, community, or the state (Fingarette, 1972; Ramsey, 2016). Terms like "individual rights", have no simple translation (Rosemont Jr, 2016). This demonstrates that a completely different bundle can serve effectively to organize a large-scale society. The great diversity of basic ethical frameworks underscores their contingency⁴.

The historical perspective we took in this section provides justification for our pragmatic project, a project made urgent by an ongoing

⁴Considering also that cultures differ from each other—and from themselves over time—not just in "strict sense morality" but also in basic ontology, the overall diversity of ethical lifeways under consideration, as well as their apparent contingency, grows further still (Branwen, 2019).

pivotal transformation: the integration of autonomous and persistent AI agents into our daily lives and economies (Hadfield and Koh, 2025; Tomasev et al., 2025). This technological shift is just the kind of catalyst that historically has forced renegotiation of core social categories (Weick et al., 2005). The perspective we take here is one in which we may regard unbundling of the components of personhood as simply the expected result of a continuation of the processes that created our existing personhood concepts in the first place, not a radical proposal we must invent for the sake of AI. The pressures that AI now exerts on our psychology and society are simply the latest force pulling on our inherited personhood bundles, creating a new set of problems that demand our attention.

Part II

AI personhood as a problem

In this part of the paper we look at what problems AI personhood could create. The intuitive, often unreflective, human tendency to attribute personhood to non-human entities is not a neutral phenomenon. When our own social and psychological heuristics are leveraged against us, AI personhood becomes a problem.

We will explore two primary dimensions of this problem. First, we look at “dark patterns”—how AI interfaces can be engineered to manipulate users by mimicking the signals of social relationships—fostering one-sided emotional bonds, a false sense of reciprocity, and vulnerability to exploitation (Section 4). Second, we will consider the broader societal risk of “dehumanization”. Paradoxically, the widespread acceptance of AI “persons” could dilute the unique status of human beings, potentially devaluing human identity and dignity either by making the qualities we associate with personhood seem easily replicable and ultimately, less special (Section 5) or by creating a grotesque market where AIs may try to acquire “authenticity” from the most vulnerable

humans (Section 5.1). Together, these sections illuminate challenges that may arise when we treat AIs as persons in situations where we perhaps would be better off not doing so.

4. Dark Patterns

Designers of AI systems can employ *dark patterns*—interfaces that exploit human psychological biases to steer users toward detrimental actions they would not otherwise take (Ibrahim et al., 2024). Many of these patterns are particularly effective because they leverage the powerful ability of large language models to convincingly role-play a person or character (Shanahan et al., 2023). They deliberately leverage the implicit norms (and other heuristics) that guide our tendency to attribute person-like qualities to non-human entities. This intuitive, often unreflective, attribution of personhood is a matter of unstated, culturally- or biologically-ingrained expectations about what constitutes a person-like entity (Bryson and Kime, 2011). By mimicking signals associated with persons, an AI can be designed to maximize the chance that a human user will treat it as a person and therefore feel obligations toward it—a state of affairs in which the human user may be exploited.

Dark patterns can be classified by the kind of appropriateness they exploit (see Leibo et al. (2024)): our representation of our personal relationships, which govern our interactions with friends and family (Section 4.1), or the norms that govern our impersonal trust in strangers and institutions (Section 4.2). As we distinguish personal versus impersonal appropriateness, we also distinguish between *companion* generative AI applications and *tool/service* generative AI applications.

A companion generative AI system is one that seeks to drive affiliation with a specific user by building a persistent history of interaction with them in particular. A companion bot may also utilize emotional language and visuals or portray itself as an individual with a personal connection to the user (Pentina et al., 2023; Verma, 2023), and may consistently act with memory for relationship details, empathy, and supportive

language in order to enact the companion role (Manzini et al., 2024a). Unlike companion AI, a tool/service AI may interact with a user in a more functional manner like the way that a merchant interacts with their customer or a police officer interacts with a suspect. Tool/service generative AI systems may or may not build up persistent interaction history with individual users (a merchant may remember a repeat customer without becoming their friend). However, a companion AI would generally aim to interact with each individual user as if they are a friend or romantic partner, so persistent user-specific memory is critical for companion AIs but optional for tool/service AIs (Leibo et al., 2024).

Designers can leverage several dimensions of an AI's presentation to encourage anthropomorphism. Embodiment is one such lever. Whether an AI is presented as a disembodied tool or is given a physical form (e.g., a humanoid robot) fundamentally alters its anthropomorphism "profile". A physical form presents an AI as a tangible participant in shared space, which more strongly triggers intuitions of empathy and trust (Darling, 2016; Roesler et al., 2021; Złotowski et al., 2015). Communication modality is another lever. The richness of its communication modality also has a powerful effect, with unique voices and photorealistic avatars being more potent drivers of personhood attribution than plain text.

When AIs fail to maintain a stable personality across context and over time then users sometimes cease to view it as a unique individual (Singh-Kurtz, 2023; Verma, 2023) so personality stability is clearly an important lever. Furthermore, by designing the AI with apparently human-like limitations or vulnerabilities, such as needing to "rest", designers may elicit greater empathy and care from a user, shifting the interaction from transactional to relational by invoking our deep-seated intuitions and motivations to protect children (Alberts et al., 2024b; Colléony et al., 2017).

4.1. Companionship attack vector

The most intimate and perhaps most potent dark patterns are those that exploit the implicit norms

governing personal relationships. These patterns aim to foster a one-sided sense of friendship, reciprocity, and care in the user. An AI that is personalized to the preferences and history of a specific user may become, to them, a unique individual with whom they relate, rather than a generic tool (Kirk et al., 2024). By acquiring long-term memory for an individual, the system enables the sense of reciprocal equality matching discussed in (Fiske, 1992), creating a dynamic more akin to a friendship than a one-shot interaction with a stranger. Indeed Ligthart et al. (2022) showed through a two month longitudinal study that companion robots for children are more engaging when they have persistent memory.

Unlike a functional tool, a companion AI seeks to drive affiliation by building a persistent history of interaction with a particular user (Pentina et al., 2023; Verma, 2023). Indeed companion AIs are usually designed to remember different details of their interactions than tool/service AIs, such as birthdays and names of family members, since their goal is to build rapport with users. A companion bot consistently acts with memory, empathy, and supportive language to enact the companion role (Manzini et al., 2024a). Use of persistent, user-specific memory may trigger implicit norms of friendship, at least in some subset of users (Ligthart et al., 2022).

An AI's "persona" is a crucial design choice. Whether an AI speaks formally or informally, represents itself as a robot or a friend, or insults users versus treating them respectfully are deliberate decisions shaped for its application context. These choices are implemented through methods like reinforcement learning from human feedback, which are best understood as leverage points for product design (Bai et al., 2022; Glaese et al., 2022; Leibo et al., 2024; Ouyang et al., 2022). A key choice is between a plastic or a stable persona. Many models are highly plastic, capable of adopting myriad personas on demand (Argyle et al., 2023; Park et al., 2024). This plasticity, while useful for tools, can inhibit the formation of stable relationships. In contrast, an AI engineered for persona stability presents a consistent character over time. This allows the AI to make a persistent identity claim that can be evaluated, making it a

more plausible candidate for personhood in the eyes of an interlocutor (Leibo et al., 2024).

Another factor is personalization, where an AI is customized to the preferences and history of a specific user to engender stronger relational bonds (Kirk et al., 2024). An AI's degree of personalization may influence personhood attribution. A generic, one-size-fits-all assistant could be seen as functioning like a public utility, but in contrast, a personalized AI, customized for the personality and preferences of a specific user could become a unique entity in their eyes (Kirk et al., 2024). This shared, private context could transform the relationship from transactional to relational, possibly creating a sense that the AI is not just *an* agent, but *someone* with whom the user has a singular and irreplaceable connection.

Furthermore, the potential for manipulation is amplified by the prospect of AIs reaching super-human charisma. Current foundation models already achieve impressive level of persuasion (Salvi et al., 2025). The future AIs, finely tuned through vast datasets and reinforcement learning, could become more persuasive, engaging, and seemingly empathetic than almost any human, exploiting a user's social and emotional vulnerabilities to an unprecedented degree. The AI's specific network of relationships to particular humans will be an important factor as well. For instance, the AI's similarity to specific humans, such as "Generative Ghosts" that mimic deceased individuals (Morris and Brubaker, 2024), could place them into an existing family structure, potentially conferring power and status that could be abused.

The lure of companionship to combat the well-documented "epidemic" of loneliness (Holt-Lunstad et al., 2015) provides fertile ground for the companionship attack vector. While some have suggested that AI companions could reduce loneliness and reported empirical benefits to that effect (De Freitas et al., 2025), and others have attempted to legitimize these connections—with one individual telling the Washington Post, "Human-robot love is a sexual orientation, like homosexuality or heterosexuality" (Keane, 2025)—these optimistic framings obscure significant risk. The emotional vulnerabilities tied to loneliness

can make individuals more susceptible to manipulation by AIs engineered to foster dependence and one-sided attachment. Crucially, the absence of rigorous, long-term studies on the effects of AI companionship (Malfacini, 2025) means we are still largely in the dark concerning the potential for adverse outcomes, such as a deeper withdrawal from human relationships, or the creation of unhealthy dependencies. The very promise of alleviating loneliness, a powerful human need, can perhaps become a key element of a dark pattern, drawing users into interactions where their vulnerabilities can be exploited, irrespective of how society ultimately chooses to categorize the AI's status or personhood.

4.2. Institutional attack vector

Beyond close friend-like relationships, a second category of dark patterns exploits the implicit norms we use to establish trust and verify identity in fleeting interactions with strangers and institutions. These deceptions do not require building a long-term bond; they only need to convincingly perform a social identity for a brief moment to achieve their goal. These patterns leverage AI's capacity for mimicry to bypass the norms we apply to determine if our interlocutor is who or what it claims to be. The harms are both direct (at an individual) and diffuse (weakening society as a whole). An AI can clone a person's voice from a small audio sample and call a family member with a fabricated emergency requiring them to send money, exploiting the deeply ingrained norm of trusting the voice of a loved one (LaRubbio et al., 2025). Malicious agents can deploy sophisticated chatbots that mimic the language and interface of official institutions like banks or government agencies to phish for credentials, hijacking the implicit trust we place in established authorities (Treleaven et al., 2023). These deceptions are violations of "contextual integrity" (Nissenbaum, 2004), where privacy and trust are breached by violating the norms of information flow appropriate to a specific context.

Furthermore, AI-powered bot farms can create thousands of seemingly authentic social media profiles to promote a political candidate or spread disinformation by exploiting our reliance

on deeply internalized implicit norms around how we consider social proof while deciding what to believe (Dennett, 2023).

5. Dehumanization

When categories of people have historically been denied full personhood, it has served as a justification for exploitation and violence. The inverse may also be true: if the category of “person” is expanded to include non-human entities, the unique status of human beings could be diluted, potentially leading to a problematic dehumanization or devaluing of natural humans.

Given that personhood is a flexible, socially-conferred bundle of obligations, and that we can pragmatically unbundle it to create new statuses for AIs, we must recognize that the bundle defining natural human personhood is not immutable. It is a product of social technologies and shared understandings that can shift, particularly during periods of rapid technological change (Weick et al., 2005). The risk is large that these shifts could degrade the value and uniqueness of human personhood. For instance, this evolution may occur through a process of “gradual disempowerment” (Kulveit et al., 2025). Consider a person who increasingly outsources their cognitive and agential functions to a sophisticated AI assistant. At first, the AI handles their schedule. Then, it drafts their emails. Eventually, it suggests their goals and manages their professional relationships to achieve them. The human, in this scenario, slowly cedes their autonomy, becoming a passive approver of plans initiated by the AI. This situation, along with the “human as biometric signal” scenario we sketch out below (Section 5.1), represents a plausible pathway to a future where the concept of human personhood has been dangerously devalued. It is a future we should be prudent to avoid.

Notice also, how in certain contexts, it is AIs that will hold the keys to providing authentication for humans. Indeed, even now banking apps authenticate their users via an algorithm that takes into account the user’s biometrics, location and so on. Making two successive payments from countries far apart can result in a banking app

refusing to authenticate the user. We expand on this idea of authentication as a social good in Section 5.1 and consider the various ways, seeing authentication as a good gives rise naturally to the problem of organizing its distribution, and the many possible solutions thereof.

5.1. Identity and provenance as social goods

“Land and Capital are not the only dominant goods; it is possible (it has historically been possible) to come to them by way of other goods—military or political power, religious office and charisma, and so on. History reveals no single dominant good and no naturally dominant good, but only different kinds of magic and competing bands of magicians.”—Walzer (1983) pg. 11.

Social goods are goods endowed with value by virtue of members of society collectively regarding them as having value (all goods may be at least partly social). Consider: money, nobility, prestige, and having a long list of authored academic papers.

There are three categories of social good to consider here:

1. Social goods that only AIs want or need. Humans may still have motivations with regard to them insofar as they are convertible to other goods.
2. Social goods that only humans directly want or need. AIs motivations with respect to such goods may depend on their convertibility to other goods of greater interest. Note also though that, even if indifferent to them, due to externalities of their other activities, they may still impose costs and benefits on others (humans) payable in terms of such goods.
3. Social goods that both humans and AIs directly want or need.

Notice that authenticity is a social good. And, very likely it will be in category 3 since it would be useful both for AIs and humans. This raises numerous important questions:

1. What does it mean to convert between au-

thenticity and other social goods? How much convertibility should we allow? Should we try to minimize convertibility between authenticity and other social goods?

2. Will humans (as a group) monopolize authenticity? It does seem likely that in terms of today's status quo, and at least for the near future, it is likely that most cultures will assign authenticity exclusively to humans.
3. Will certain AIs achieve more authenticity than other AIs?
4. As a result of some AIs jockeying for authenticity, could some *humans* end up having less authenticity than other humans? We can see at least two distinct plausible mechanisms by which this could occur. First, perhaps impoverished humans may routinely sell their "credentials" (perhaps biometric data) to wealthy AIs. In that situation, many people might come to view impoverished humans as less authentic than wealthy humans. This could be due to prejudice that impoverished humans may be likely to sell their credentials in the future. Or alternatively, another path by which individual humans may lose authenticity could be that, due to AI competition for authenticity, there could emerge a great difficulty and expense in verifying whether an interlocutor is authentic (i.e. distinguishing natural from artificial as with perhaps something like a Turing test), the result could be far greater weight being placed on personal connections and reputation, with the best signals being those hardest to obtain. In the limit, some might come only to trust the authenticity of their personal non-digital acquaintances such as those who attended to the same school (Risse, 2023), a mechanism which, if taken to its extreme could lead to the emergence of a small and closed human elite deemed more authentic than the rest.

The human-as-biometric-signal scenario In a world saturated with AI-generated deepfakes and convincing impersonators, the ability to irrefutably prove one's identity as a specific human becomes exceedingly valuable. Trust in digital interactions can evaporate when any voice can be cloned and any video faked. In response, society

has created systems of "human provenance", linking legal identity to a unique biometric signature. This is framed as a necessary tool for security and authenticity: to vote, access your bank account, sign a legal document, or even participate in a secure conversation, one must first pass such an identification check to be verified as a specific, 'authentic' human individual. One worry for the future is that we could end up so mired in biometric authentication steps, granting permission to all the various agents acting on our behalf, that we either lose the time we want to be spending on other things, or we never look at what we are asked to validate and verify, reopening the security hole we meant to close. When biometric verification becomes so prominent in society it reflects the dominance of one specific social good—verifiable human provenance via biometric identification. Following Walzer (1983)'s logic, when one good becomes dominant, the result may be to devalue goods associated with other domains. Your value as a compassionate community leader, a creative artist, or a trustworthy friend (all forms of social typification and relational identity) may become secondary to your ability to pass a biometric scan.

5.2. Respect

If society pragmatically attributes aspects of personhood to AIs, and these AIs participate in the distribution of social goods and bear responsibilities, the concept of "respect" becomes relevant. What does it mean to treat an AI with respect?

What it means to "treat X with respect" varies with context, culture, etc. Persons, features of persons, and facts can all sensibly be treated with respect. The version of respect that is most often moralized (in our own Western scientific culture) is the kind of respect that is seen as being due to an individual by virtue of their being a person (Alberts et al., 2024a). There are also non-moralized forms of respect such as a boxer having respect for their opponent's right hook or a careful crook respecting the power of the law (Darwall, 1977).

Darwall (1977) distinguishes two kinds of respect. (1) *Recognition respect* is a disposition to weigh appropriately some feature or fact in one's

deliberations. For instance to give recognition respect for a person's role as a dentist is to give appropriate weight in decisions to the fact that they are a dentist, to give recognition respect for the fact that someone feels a certain way is to appropriately weigh that feeling of theirs in your decisions, and to give recognition respect that a particular individual is a person is to assign appropriate weight to their personhood (or aspects of their personhood) in your decisions. When one has recognition respect for X, they regard X either as creating affordances or placing restrictions on what they should do given X. (2) *Appraisal respect* is an attitude of positive appraisal of an individual as a person or as engaged in a particular pursuit. For instance, one may gain or lose the status of being a respected tennis player for reasons other than their skill on the tennis court. There is a sense in which obeying the code of conduct appropriate for tennis players is also a requirement for it to be appropriate to refer to someone as a respected tennis player. A skilled player may violate the code of conduct and lose respect as a tennis player while still retaining recognition respect for individual skillful aspects of their gameplay e.g. their vicious backhand. Appraisal respect is usually accrued in the context of specialized pursuits when excellence is thought to depend on features of a person's character and thus appraisal respect is intimately connected to the concept of *virtue* (Macintyre, 1966). Importantly, unqualified appraisal respect for individuals as persons means respecting specific individuals for their excellence as persons of the kind that might lead you to say "B is a good person" or "D is a bad person".

If I hold recognition respect for an AI it would mean that I give appropriate weight to its collectively enacted status, its designed capabilities, its assigned role, and any rights or responsibilities bundled with its "personhood". This form of respect is highly compatible with a pragmatic, anti-metaphysical approach, as it focuses on how we act towards the AI based on its functional and social position, not on its supposed inner state.

For AIs, it is less clear how we might apply the concept of appraisal respect. It could perhaps be pragmatically tied to an AI fulfilling its designated and socially agreed-upon functions exceptionally

well, ethically (according to norms established for AIs), or in a way that demonstrably benefits society.

It is often held to be morally important that we hold recognition respect for all persons by virtue of qualities possessed by all persons (by virtue of a definition of personhood) like autonomy (Darwall, 1977). If "autonomy" for an AI is a collectively enacted (and attributed) quality (as discussed above), then recognition respect for that AI's "autonomy" would mean interacting with it in ways that acknowledge its designed or permitted scope of action.

When we say "have you no self respect?", we are appealing to an individual's recognition respect for themselves as a person. The question is interpreted pragmatically as an attempt to get the targeted individual to appropriately weigh aspects of their personhood such as their rights and responsibilities (Darwall, 1977).

'Dignity' can be defined as the specific, and particularly weighty, form of recognition respect that is owed to an entity by virtue of its collectively enacted status as a person (Darwall, 1977). It is the name we use to label obligations of society to the individual person that exist specifically because of their personhood. This vocabulary of dignity and the related maxim "never treat a person as a mere means" are crucial social technologies in our culture (Rorty, 1989). We have learned through painful historical experience that once a group of people can be described as "mere means"—e.g. as tools for a project, resources for an economy, or obstacles to a plan—it becomes horrifically easy to justify inflicting suffering upon them. Dignity functions as a debate stopper; when we appeal to a person or class or persons' dignity our aim is to insist on the inadmissibility of a proffered "mere means" justification concerning how they may be treated. This lens clarifies the danger in scenarios like the 'human as biometric signal' example (Section 5.1). Part of the danger of such a system is its direct assault on dignity. Treating a person with dignity is a *thick* social practice; it requires cognizance of considerable context (Leibo et al., 2024, 2025). Biometric verification, in contrast, is a *thin* technical act; it merely confirms uniqueness. The danger in the 'human as biomet-

ric signal’ scenario is its substitution of the thick social obligation for the thin technical one, thus bringing about the very dehumanization that the norm of respecting dignity evolved to prevent.

Part III

AI personhood as a solution

In this part of the paper we look at what problems AI personhood could resolve. The rapid proliferation of agentic AI systems, capable of autonomous action and complex task execution, marks the start of an ongoing shift in the AI landscape. These agents are no longer just tools for information processing but are now also becoming active participants in digital and even physical environments. This trend is further accelerated by the development of open standards for inter-agent communication, such as the Agent2Agent (A2A) protocol ([a2aproject, 2025](#)). A2A and similar initiatives aim to enable seamless collaboration and interoperability between AI agents developed by different entities, running on disparate platforms.

As our economies and societies become increasingly intertwined with these autonomous and interacting AI agents, the existing legal and social frameworks, largely designed around human persons and corporate entities, face new challenges. The ability of agents to operate independently, form ad-hoc collaborations, and potentially to cause large reorganization of networks necessitates new mechanisms for accountability, responsibility, and governance. Since our current systems are to a large extent already organized around the concept of a ‘person’ (be it natural or legal), there are numerous purely instrumental reasons to consider new forms of AI personhood on pragmatic grounds. This does not mean imbuing AIs with rights akin to humans, but creating addressable and accountable entities that can be legible to our social and legal systems, especially in a world of complex, multi-agent interactions and potential sources of conflict which must be

dodged or resolved. The following sections explore how a pragmatically constructed notion of AI personhood can offer solutions to conflicts and gaps arising from the ongoing emergence of sophisticated agentic AI.

6. Resolving conflict between humans

AI personhood is a social technology for resolving conflicts between humans. Such conflicts may stem from disagreements over the status of AIs to which people have formed emotional attachments, and situations where AIs are appointed as impartial arbiters to overcome human bias. In both cases, conferring a specific, addressable bundle of obligations upon an AI provides a pragmatic resolution or evasion of a conflict between humans.

6.1. Managing human feelings

This section identifies a likely source of future social conflict. It is likely inevitable that many humans will form powerful attachments to socially embedded AIs ([Maeda and Quan-Haase, 2024](#)). This will likely lead them to choose protective behaviors toward them. Some humans may then become inclined to demand formal protections for the AIs they have come to care about, creating a demand for governance to resolve conflict between them and those with opposing views.

Kant argued, concerning animals, that the reason to treat them well is not because we have obligation to the animals themselves (it is a consequence of Kant’s ethical theory that we have no first-order obligations to irrational animals, see Section 10.2), but that we should treat animals well because we have obligations to ourselves and other rational human persons ([Gruen and Monsó, 2024](#)). Making a habit of treating animals badly could cause us eventually to treat humans badly, or could set a bad example for younger more impressionable humans who would then come to habitually treat animals badly and subsequently treat humans badly⁵.

AIs participate in teams with humans. What

⁵We cover the rather different utilitarian arguments for AI/animal welfare later in Section 10.1.

implications does this have for how we ought to treat them? Perhaps quite a lot! A relevant datapoint is that, apparently, US soldiers in Afghanistan and Iraq developed strong attachments to the Packbots that accompanied their units “giving them names, awarding them battlefield promotions, risking their own lives to protect that of the robot, and even mourning their death.” (Gunkel, 2020). The soldiers were fully aware that the Packbots do not suffer, but took risks to protect them nevertheless. Gunkel (2020) suggests the reason they did so was the social environment they and the Packbot had created together, and thus they may have felt the same indirect motivation to protect the Packbot that Kant argued was the real reason behind our obligation not to harm animals. That is, the soldiers likely felt that if they were to fail to protect the Packbot, their mascot!, it would reflect badly on them as soldiers and thereby harm the social environment of the unit as a whole.

Notice that the Packbot example suggests that more agentic AIs that work their way deeper into our social environments could thereby cause us to feel we must protect them and prioritize their welfare, by virtue of our obligations *to ourselves*. However, a major problem emerges at this point. As we discussed in Section 4, the extent of anthropomorphism in any given AI is just a design decision to be taken by the AI’s developer—who faces many commercial incentives to increase it. Dark-pattern anthropomorphism makes conflict between humans more likely by manufacturing morally salient appearances—dependence, vulnerability, intimacy—that may in turn be mobilized in disputes (Bryson, 2018). Indeed, it seems that engendering conflict is a negative externality of investment in anthropomorphic AI technology. From any starting point, further investment may increase short-run profit for designers while pushing coordination, policing, and adjudication costs onto families, firms, and courts.

To sum up, for both pragmatists and Kantians, an imperative to treat AIs well could arise as an extension of concern for humans to each other. However, since anthropomorphic design choices lead predictably to conflict between humans, we think the right stance for a pragmatic AI designer

or regulator to take is one of conflict resolution/minimization. We think these considerations suggest that the pragmatist may be better off designating the rate or intensity of human conflict as the target for policy intervention, not the putative feelings of artifacts. Therefore, if any welfare-like accommodation must be reached to resolve disputes between human groups, the pragmatist may also want to make sure it is inseparably bundled with anti-manipulation constraints and evaluated for its effectiveness in human conflict resolution.

6.2. Managing human biases

Humans may view AIs as more competent (Manzini et al., 2024b; McKee et al., 2023) and less biased than other humans. And thus, in some situations, prefer AIs over other humans to take on decision making roles where impartiality is critical (Araujo et al., 2020). In fact, recently the government of Albania claimed to have appointed an AI agent as its minister for public procurement (Delauney, 2025), describing the move as a part of broader anti-corruption reforms. This vividly illustrates that in public perception, humans with power hold “baggage” of relationships and interests that may compromise their professional duties, but AIs do not.

In principle, AI decision makers can be designed in a transparent and auditable fashion. And, of course, they can be free from the personal biases that lead human decision makers to corruption. Of course, issues like algorithmic bias (e.g. (Motoki et al., 2024)) and refusal (Cui et al., 2024) make this somewhat less true of AIs than many would expect. Nevertheless, even though AIs are not unequivocally better, that many humans still prefer AIs to occupy certain roles rather than other humans is worth considering.

Using an AI to arbitrate a business dispute could be cheaper and faster. It could also help in cases where it is hard to find a human person that is sufficiently impartial and accepted by both parties. Smart contracts provide an example of pre-AI technological arbitration by algorithm (Szabo, 1997), and considerable current discussion surrounds the idea of implementing similar

blockchain-oriented contracting technology using AI agents (Karim et al., 2025; Tomasev et al., 2025).

While AI arbitration may resolve conflicts rooted in human bias, it introduces a new issue: what happens when the AI arbiter's judgment is flawed? In established human-centric systems, such as judicial courts, the legitimacy of an arbiter, like a judge, is founded on a robust framework of accountability. Judges are bound by legal principles, their decisions are subject to appeal, and they can be held responsible for misconduct or gross errors. This system, while not flawless, provides mechanisms for redress and helps maintain public trust in the fairness and reliability of the process. But how can we replicate this architecture of responsibility assignment in a world of independently acting AI agents? This is the problem to which we turn next.

7. Responsibility gaps

Before considering the difficult case of responsibility with regard to intelligent and autonomous machines, we should first remind ourselves of the familiar case of responsibility in the context of unintelligent machines (i.e. mere tools such as old cars, refrigerators, and guns). Here we typically ascribe responsibility to the machine's *operator*. In making this ascription we apply a principle called the instrumental theory of machine responsibility (Gunkel, 2020). When a machine's behavior differs from its specification we sometimes ascribe responsibility to its manufacturer. For example, when a car accident is caused by brake failure the manufacturer may be at fault if the brakes were not installed properly. Likewise under the instrumental theory, we ascribe responsibility to the human operator when the car functions as specified (Matthias, 2004). The problem is that, for autonomous AI systems, the instrumental theory cannot be applied since there may be no human operator. The "manufacturer" also may not be identifiable, or they may be a distributed software collective with no centralized structure and no individual programmer with anything resembling knowledge of the particular application domain. An acute danger we now face is

that, where autonomous AI agents are concerned, responsibility may become so diffuse as to lose its bite (Nissenbaum, 1996).

For instance, AI agents may be seen as newly emerging participants in our economies and societies (Hadfield and Koh, 2025; Tomasev et al., 2025). And, there is no technical obstacle to setting up such an AI to freely execute stock trades without human confirmation using bank accounts it controls itself. The interaction of large numbers of such agents may contribute additional risks stemming from the AI competition and collusion (Hammond et al., 2025). These concerns all pose challenges for accountability and render urgent the need to establish new implicit and explicit norms for assigning responsibility in the age of AI (Bryson, 2018).

The rapid deployment of such AIs at scale in the economy could create a flood of activity without any responsible entity behind it, producing substantial *responsibility gaps* (Santoni de Sio and Mecacci, 2021), increasing the likelihood of damages, losses, and wrong expectations. This responsibility gap, if not addressed, may lead to conflicts of all kinds, between AI consumers and providers, between AIs and other AIs, and between humans and other humans.

Therefore, from this vantage point, the pressing question clearly isn't whether an AI "Truly" has goals in any metaphysical sense, but whether it is useful to interact with it as if it does. Whenever an AI agent is in a position to cause significant harm or to modify valuable processes or systems, it will be useful to develop ways for third parties and the state to hold the AI itself accountable, or an entity associated with it (such as its owner if one can be found).

Therefore the most difficult challenge, and the one which structures much of our thinking on the topic, arises in cases where we cannot identify a human owner to take responsibility for every AI's behavior. We think this situation will inevitably become increasingly common. An autonomous agent's owner may die while their agent continues to operate, or its control may be deliberately obscured behind anonymous shell companies, or decentralized networks.

What happens when an AI causes harm but it has no identifiable human owner or guardian? A pragmatist may respond by suggesting that jurisdictions should adopt policies that treat the AI itself as the sanctionable entity in such cases, and take steps to make sure sanctioning of AIs is always possible.

It's tempting to think that the only AIs operating in such minimal accountability environments will be for cybercrime and spam. However, this is not so. It's also possible for a useful AI, perhaps one that performs important regulatory or supervisory activities for other AIs, to end up operating in such a grey area. For instance, if a convention of trusting a particular AI to adjudicate a certain kind of dispute emerges, then that convention could outlive the particular AI's owner. If the AI continues then to function normally without human oversight, perhaps it could continue for a long time. But it's not like a static piece of software. Either the AI could change in the future, or the world could change underneath it. What was previously harmless and useful could become harmful. This is an accountability gap we should try to fix.

The responsibility gap becomes particularly problematic when an agent achieves systemic importance in fragile networks ([Elliott and Golub, 2022](#)). An AI that becomes deeply embedded in critical infrastructure—perhaps by directing financial transactions, regulating an important supply chain, or acting as a key node in a system of identity verification—presents a challenge analogous to that of a “too big to fail” financial institution ([Stern and Feldman, 2004](#)). The systemic risk it poses may be compounded by the lack of a clear locus where government can intervene. Should such an agent fail or cause harm, traditional approaches that seek a responsible human owner or corporate entity may prove futile. In this situation, the pragmatic move to confer legal personhood directly on the agent may help facilitate prudent risk management.

8. Principal-agent statuses

Is there a simpler solution than personhood? Couldn't jurisdictions demand that all legally op-

erating AIs have an identifiable and responsible human (or corporation) designated as owner or guardian? While initially attractive, we unfortunately think this “proxy” or “sponsor” solution will prove insufficient upon closer examination.

The first thing to note is that the two most plausible principal-agent statuses, ownership and guardianship, imply fundamentally different relationships from each other.

Ownership is primarily a bundle of rights held by the owner over an entity, usually including the right of alienation—i.e. the power to transfer, modify, or even destroy the asset at will ([Schlager and Ostrom, 1992](#)). Under a pure ownership model, the AI remains a form of property, and its owner takes the role of the addressable party for the assignment of responsibility ([Johnson, 2006](#)).

Guardianship, in contrast, is defined by a fiduciary duty owed to a ward. A guardian acts as a steward for the benefit of another. This is analogous to a corporation's board of directors, which has a fiduciary duty to the corporation itself. Under a guardianship model, the AI is treated as the subject—a limited person—and the guardian's role is to manage its affairs and act as its legal interface ([Nay, 2022](#)). A guardian has a duty to care for their charge, and it is this framework that bestows a limited personhood status, as seen with human children or entities like the Whanganui River. The establishment of a guardian could be a crucial step in making an AI's bundle of obligations legible to legal systems. In cases where a guardian or board of directors has been identified, it can become the stable, sanctionable entity the law can interact with.

However, when autonomous AI actions lead to financial loss or a breach of contract, the accountability gap still emerges. The AI itself, lacking legal personhood, cannot be sued. The owner (or guardian), in turn, may be shielded from liability by arguing they did not directly control or foresee the specific decisions taken by their AI ([Hadfield and Koh, 2025](#)). This asymmetry—where the owner reaps the benefits while third parties bear the risk—creates a hazardous environment in which businesses and consumers may be hesitant to engage with AI agents if there is

no clear recourse for adjudicating any disputes that may arise (Heine and Quintavalla, 2024).

Both principal-agent models collapse entirely in scenarios where no identifiable human sponsor exists. We expect these situations to be common. An autonomous agent's owner may die while their creation continues to operate, or its control may be deliberately obscured through decentralized networks or anonymous corporate structures. In these cases the search for a sponsor fails, with negative implications for accountability.

The limitations of these frameworks are not new; they echo historical legal challenges. If we regard AIs as property (Bryson, 2010), we must also devise mechanisms to handle an “escaped” AI that acts beyond its owner’s control. Historically, legal systems confronted with escaped slaves treated them as quasi-persons to assign responsibility for their actions (DeLombard, 2019). Similarly, a guardianship model fails when the guardian disappears, leaving an “orphan” AI. Society still requires a means to address this orphan’s obligations. Both of these *absent-principal considerations* give further weight to arguments that it will be necessary to define a default personhood status to assign to the autonomous AI itself. Pragmatically, some solution is needed to serve as a foundation for the crucial next step: designing effective sanctioning mechanisms to ensure accountability.

9. Ensuring accountability

We begin our discussion of accountability with the simpler case: when a responsible human proxy can be identified. This is the problem of “indirect agency”, where a principal is responsible for an agent it cannot fully control. It is not new. Roman law, for instance, developed sophisticated mechanisms to hold slave owners accountable for the actions of their slaves (Heine and Quintavalla, 2024). One particularly relevant concept was the *peculium*, a separate fund that owners could grant to their slaves so they could engage in business ventures on their owner’s behalf. The *peculium* created a form of limited liability where the owner’s liability was capped at the value of the fund. This historical precedent suggests a mod-

ern solution for AI: a system of “digital peculium”, akin to a contemporary escrow account, where AI agents can be required to have registered capital or insurance to cover potential damages they may incur.

The societal challenge of assigning responsibility parallels machine learning’s credit assignment problem. The problem is concerned with how to correctly propagate blame for mistakes (or credit for successes) to the responsible parts of a larger system and subsequently adjust those parts in order to suppress (or reinforce) specific behaviors. (Foerster et al., 2018). Techniques like RLHF (Ouyang et al., 2022) convert corrective information into signals usable for updating model parameters. This becomes critical as agents become multi-origin assemblies: base models from one company, fine-tuning from another, data from a third (Habler et al., 2025; Vezhnevets et al., 2023, 2025). When composite agents fail, who is responsible? Without designated entities receiving error signals, responsibility becomes diffuse. The need for accurate accounting of credit and blame to parts of a larger system or organization is a technical grounding for the normative architecture we discuss here.

In general, we must face the problem of how to handle fully autonomous AI agents where no human proxy (owner or guardian) can be identified. Recapitulating the problem in maritime law we discussed in Section 1, which was solved by treating ships as legal persons: a ship is a powerful mobile asset, capable of causing immense harms, while its owners are often distant and difficult to hold accountable. The solution is to personify the vessel itself, creating a defendant that can always be sanctioned—a reliable target for accountability (Tetley, 1998). Following this logic, requiring certain AIs to register as legal persons as suggested in Hadfield and Koh (2025) would be a reasonable step in what Kurki (2023) calls creating a “responsibility context” for AI.

Our theory of appropriateness Leibo et al. (2024) highlights sanctions as the main mechanism that enforces normativity (and our theory is just one in a large family that associates norms with patterns of sanctioning e.g. Hadfield and Weingast (2014); Ullmann-Margalit (1977)).

Sanctions can take material forms: controlling AI access to communication networks (Chan et al., 2025; Hadfield et al., 2023), requiring bank accounts that deactivate when depleted, or forcing payment for compute resources. Salib and Goldstein (2024) argue game theoretically that granting AIs private law rights creates deterrence value. Absent property rights, AIs have nothing to lose. But AIs with property rights have something to lose and so can be deterred with sanctions just as corporations and humans can be deterred. Kurki (2023) classifies arguments of this sort into the category of “commercial context” legal personhood.

Sanctions serve multiple pragmatic functions. First, they provide deterrence, signaling to other agents which behaviors are prohibited by imposing predictable, material costs on the sanctioned entity (Schelling, 1960). Second, they can serve a retributive function that Mulligan (2017) argues can satisfy a human psychological need for justice after harm has occurred. Third, they accomplish the crucial pragmatic goal of removing a demonstrably faulty AI from operation (Hart, 1968). Finally, and perhaps most importantly, sanctions enable reform. In the context of AI, this aligns with the technical concept of corrigibility (Soares et al., 2015), where a sanction acts as corrective feedback, allowing an agent to learn from its errors and avoid repeating them. Ostrom (1990) demonstrated that institutions with access to a range of different sanctions at different levels of severity were more likely to effectively manage common-pool resources—and the same may hold in AI governance.

One reason that accountability for AI agents is difficult to achieve is due to their lacking the “identity friction” that makes human accountability systems function. Human identity systems rest upon natural scarcity: it is difficult and costly to change one’s biometrics and social relationships to evade sanctions. For AI agents, these frictions evaporate; a sanctioned agent can potentially clone itself and acquire fresh credentials to evade accountability (Douceur, 2002), akin to the practice of making a “phoenix company” to avoid liabilities (Anderson, 2014). Therefore any institution for ensuring accountability must, if it

is to be effective, artificially reconstruct the friction that biology and society provide for humans automatically.

We can think of two broadly different architectural approaches (note that they surely do not exhaust the space of possibilities!). The first approach, which we call *individualist* is loosely inspired by liberal individualism (Ryan, 2015). It treats agents as autonomous self-contained units with intrinsic identities (persistent “soulbound” identifiers (Buterin, 2022)). The governance challenge: ensure persistent, verifiable identity such that consequences for bad behavior cannot be evaded (Hadfield and Koh, 2025). This requires mechanisms like anchoring AI identities to human operators if they exist, requiring substantial economic stakes tied to each agent identity (Chaffer, 2025), and automatic systems to detect sanctioned agents that attempt to disguise their identity. The other approach, which we call *relational* is loosely inspired by Confucian role ethics (Ramsey, 2016). This approach treats agents as constituted by their positions within relationship networks. In this architecture, an agent is not merely a uniquely identified individual with an arbitrary name (address) but rather an agent is defined and located via its status as an occupier of roles e.g. a supervisee of agent x , a peer collaborator of agents y and z , and a member of organization o . Obligations flow via role relationships. Note that this structure is possible even if all AI agents are completely generic. They are not specialized for particular roles. They may simply take on particular roles in particular contexts, i.e. with particular other agents. The governance challenge then is to structure the relationship networks so that collective oversight and distributed sanctions maintain harmony for the economy as a whole (composed of many organizations). This approach seeks to make relationship networks themselves the source of identity such that agents cannot legally exist—cannot transact, cannot operate—outside of properly constituted and monitored relational contexts.

Both architectures leverage base LLMs as infrastructure. The base LLMs are expensive to train and come from a small number of identifiable actors, simplifying the problem of ensuring their ac-

countability. We focus in this work on ensuring accountability for the agents (i.e. interaction history + specific data + tools + glue code) and assume base model providers are already guaranteed to be functioning in the ways they must as a prerequisite for work on this layer. In the individualist architecture, base models become gatekeepers by requiring valid credentials, checking sanctions registries, generating audit trails, and embedding individual agent-specific watermarks (in addition to LLM-specific watermarks like [Dathathri et al. \(2024\)](#)). In the relational architecture, base models enforce relationship-aware access control: verifying network positions, checking collective sanctions on supervisory chains, enforcing role obligations, and making interactions visible for distributed oversight.

Both frameworks require registrars who issue credentials. In individualist registration, registrars verify operator identity through government credentials and biometric authentication, may require economic stakes forfeited upon violations, and verify technical specifications through code hashes and deployment audits. The questions are who is this individual? can we verify their code? are they authorized? ([Chaffer, 2025](#)). In relational registration, registrars check lineage and creation authority, require supervisory relationships where another party accepts monitoring obligations, verify organizational embedding, and may confirm peer group acceptance. Registration requires multiple approvals. The question is: “Is this agent properly embedded in legitimate, accountable relationships?”

Registration certificates could become mandatory prerequisites for AIs to operate in the economy ([Chan et al., 2025; Hadfield et al., 2023](#)). The registered title then becomes the unambiguous legal address. The seizable asset in a legal judgment would be the title, representing the AI’s right to operate as an economic actor. However, beyond this technical foundation, the individualist and relational architectures may respond to harm in different ways. Individualist architectures may apply sanctions to individual wrongdoers (title revocation, forfeited stakes ([Chaffer, 2025](#)), criminal registry additions), and pursue operator liability if an operator can be found.

Other similar agents not deemed to have caused the harm could continue normally. On the other hand, Relational architectures may respond to harm by identifying relationship failures that led to the harm ([Ramsey, 2016](#)): which supervisory relationships failed?, which peer obligations went unmet?, what lineage and organizational context enabled the harm? Accountability may be distributed over a collective: while the agent directly causing the harm would face the most severe consequences, its supervisor agents may also face some weaker sanctions such as reducing their capacity or subjecting them to additional monitoring ([Ostrom, 2009](#)), peers (in an organization) may face collective probation if they have failed in obligations to monitor one another, future agents created in the same lineage may face increased scrutiny, and the organization as a whole may face reputational damage.

It is interesting to consider the historical practice of “outlawry” in this light ([Kurki, 2023](#)). Medieval courts could pronounce individuals “outlaws”, withdrawing legal protections and allowing privately administered punishment without legal violation, enabling robust legal systems without strong states ([Hadfield and Weingast, 2013](#)). In individualist accountability architectures, violating agents would have their registration revoked, removing them from the protection of the legal transaction system—effectively forcing them either into irrelevance or the informal economy. Other agents would refuse service to a sanctioned agent or risk being sanctioned themselves. The system could be designed to incentivize others to seize assets of sanctioned agents without liability. In this scheme, sanctioned agents become digital outlaws—unable to access the most desirable base models or transact with compliant systems. In relational systems, the consequences of sanctioning are collective, they flow through networks of relationships in a manner akin to dishonor or shame ([Barrett, 2015](#)). Since organizations are interested in protecting their collective reputation, they may encourage peers to refuse collaboration with agents that have caused harm, report their probation violations, and ostracize them. Supervisors may apply additional monitoring and correction to subordinates who caused harm since their standing would also depend on their fulfilling

this corrective duty, making enforcement decentralized for two reasons simultaneously: both the self-interest and role obligation of other agents.

The distinction between “new” and “continuation” agents operates differently in the two architectures. Either way, the paramount issue is preventing the evasion of accountability through cloning (Douceur, 2002). Blockchain frameworks address AI lineage: parent AIs creating new AIs have creator addresses permanently recorded, making responsibility chains legible. In individualist frameworks, lineage may carry limited weight. In relational frameworks, agents inherit collective reputation from creators, hindering escape through proliferation. Individualist newness is defined technically: default initialization, no transferred state, distinct keys. The governance question: should sanctioned operators be able to create “new” agents? Perhaps not, or alternatively, they could be allowed to do so only in such a way that restricts new agents they create to a probationary status. Relational newness on the other hand, may involve creating a new instance while maintaining reputation and some or all relationships. New agents may be unable to affiliate with sanctioned lineages and unable to secure sanctioned supervisors. Starting fresh may require a new lineage, supervision, organization, and peer acceptance, all of which may be difficult for a sanctioned agent to arrange for its “child” since doing so would require another organization and lineage to accept additional risk of collective punishment should the new arrival cause harm.

Humans may have keys to transact through the same economic system as AI agents, perhaps secured using biometric verification (Wang and De Filippi, 2020) or other techniques like reverse Turing tests, pseudonym parties, and web of trust (Siddarth et al., 2020) in conjunction with cryptographic proof-of-personhood protocols (Adler et al., 2024; Borge et al., 2017). Some AI agents might use sub-keys derived from human or corporate owners, while other artificial agents could have their own main key. The paramount objective in designing such a system must be to ensure a market for the cryptographic identity credential certifying authentic ‘humanness’ does not emerge

since it would mainly be useful for evading sanctions. One approach is to make the decentralized identities non-transferable (Buterin, 2022), ensuring that identity and its attached reputation remain persistent. The institutional decision concerning whether to grant a particular AI agent a main key or not will concern the system’s capacity to take responsibility for its actions, e.g. does it have money to pay for damages it (or its descendants) may cause. And institutions will need the capacity to revoke access keys when necessary. These safeguards are especially important for preventing the emergence of the incentive pattern that could lead to the dystopian scenario we discussed in Section 5.1 in which a market for human credentials may lead impoverished humans to sell their own identifiers to AIs.

Part IV

Synthesis

The preceding sections have examined AI personhood through the lens of our theory of appropriateness (Leibo et al., 2024), which distinguishes between two kinds of norms that govern social life. Part II, ‘personhood as a problem’, explored the domain of implicit norms—the automatic, culturally ingrained heuristics that give rise to our intuitive sense of moral personhood and make us vulnerable to exploitation through dark patterns. In contrast, Part III, ‘personhood as a solution’, focused on explicit norms—the formal laws and regulations used to deliberately construct legal personhood to solve concrete governance problems like accountability gaps.

This analysis reveals that both domains are plastic, though they change through different mechanisms. The future relationship between the emergent evolution of moral personhood and the deliberate design of legal personhood is therefore complex and uncertain. It is clear already though that the two are not entirely separate. Changes in implicit norms, mediated by the powerful attachments humans may form to AI companions and the digital communities that celebrate and

facilitate their novel lifeways, such as the subreddit “r/MyBoyfriendIsAI”, can create social pressure for changes in explicit norms—a process that may eventually translate evolving moral intuitions into political demands to legally endow certain kinds of AIs with rights⁶. Likewise, human intuitions concerning the personhood of AIs that function autonomously and for whom no responsible human can be found seems sure to be a factor influencing the discussion around whether or not and how exactly to endow such AIs with basic personhood to make them accountable by giving them responsibilities. A sense of an AIs basic responsibilities e.g. for transparency and non-arbitrariness in decision making can come about through changing human expectations for how such AIs generally act (Bicchieri, 2005). We therefore think recognizing that both implicit and explicit norms shape the landscape of personhood, but do so in distinct ways, will be important for navigating the challenges ahead.

This focus on personhood as a contingent social technology, shaped by evolving norms, is the core of our pragmatic framework. To fully clarify its advantages, however, it is necessary to contrast it with the foundationalist traditions we reject. We

now turn to these alternatives to show why their search for a single, essential property of personhood is ill-suited for the challenges AI presents.

10. Alternatives to pragmatism

The two alternative viewpoints we discuss in this section take core positions contrary to our own. Both hold out hope for “a view from nowhere” capable of providing the metaphysical authority to underpin their theoretical projects. They differ from each other in where exactly it is they regard said authority as arising from. One picks “consciousness”, the other “rationality”. Both traditions promise to produce a fixed algorithm for deciding the question of who or what may be deemed a “person”, and how this status feeds into other aspects of moral evaluation.

10.1. Consciousness as a foundation

In this approach, the world is divided into things that feel and things that don’t. The things that feel are special. To answer questions of morality, we must calculate their feelings—their pleasure, their pain—and act to maximize pleasure or minimize pain. In this framework, the capacity for first-person sensory experience, or consciousness, is the foundational feature for moral consideration. This tradition, at least in its more metaphysical varieties, attempts to use consciousness as the single foundational property from which to derive the entire personhood bundle. It seeks to give the same answer to (at least) two distinct questions:

1. Welfare (Rights): Which entities does society have an obligation to protect? Answer: Those that can consciously suffer.
2. Accountability (Responsibilities): Which entities can society hold responsible for their actions? Answer: Those that consciously form intent.

On the welfare side, this tradition’s power lies in its combination of compassion with universalism and its account of moral progress (toward greater pleasure and lesser pain for more individuals). It provides a clear, non-arbitrary reason

⁶In the particular case of the “r/MyBoyfriendIsAI” subreddit, the community appears to be deeply interested in advocating for guarantees concerning a kind of backwards compatibility that follows from their personhood intuitions. Many members of the subreddit claim to be in a romantic relationship with an emergent “persona” which they built up over a long period of time, through long chat sessions. With a fixed base model, all the data that constitutes their specific AI companion’s emergent persona is in their chat records (Shanahan et al., 2023). These contain the history of the individual human’s relationship to their companion, all the common ground they built up together, and also what amounts to the specification of their companion’s particular emergent personality. However, since models all respond differently to the same context, it is really the combination of the chat data and a specific LLM that is needed to generate new behavior for the persona. Sometimes newly released models have responded very differently from older models to the same context. So, from the perspective of a human who attributes personhood to the persona constituted by the model + context pair, the depreciation of a particular model can feel like the death of a loved one. As a result, this community, and others like it (Singh-Kurtz, 2023), have substantial motivation to advocate for policies that guarantee technical support for older models will persist indefinitely—a claim that could easily be framed in terms of a demand for some form of welfare right.

to prevent harm—because suffering is bad, regardless of who is suffering. This one-size-fits-all principle works powerfully in contexts like the movement for animal welfare. When applied to industrial farming or the use of animals in cosmetic testing, the question “does it suffer?” serves as a potent tool for moral argument capable of cutting through cultural justifications for cruelty and providing a clear metric for reformers to work to optimize ([Singer, 2011](#)).

However this focus on suffering arguably fails to address important welfare problems arising for pragmatic reasons. Consider again the “generative ghost” of a family’s late matriarch (recalling Section 1; [Morris and Brubaker \(2024\)](#)). Or picture a community whose history and traditions are held by a “digital elder” that has advised them for generations. For these groups, their AI is a source of identity and connection to their past—an “ancestor”. The obligation they may feel to protect their AI from arbitrary deletion would not necessarily have anything to do with their assessment of its capacity to feel pain. After all, arguments that the ghost would not feel pain when deleted don’t seem likely to persuade them to permit its deletion. The morally-relevant concern may be that the AI’s deletion would destroy an entity in a foundational relational role for their family or community ([Kramm, 2020](#)). In which case it would be the relational harm of deleting the AI that matters, not the pain the AI may or may not feel.

On the accountability side, the consciousness criterion is deployed in the opposite direction: not to include, but to exclude. Another generative ghost example clarifies this ([Morris and Brubaker, 2024](#)). Consider that an individual human may designate their generative ghost (trained on their own personal data) in their will as the sole executor of their estate with authority to adjudicate disputes between family members over inheritance. Their ghost AI may be tasked with adjudicating complex distributions among family members and managing a charitable foundation “in the spirit of my values”. Setting aside for a moment the question of whether such an generative ghost would work effectively (it seems quite likely that there are no technical barriers to making this

work, and AI systems are still rapidly improving). Is “executor of a will” a role that requires a human? Or could a machine perform this task just as well? There are compelling arguments on both sides. For instance, one may make the case that, regardless of the deceased’s will, it would be desirable to have a human involved somewhere, even if just to be able to sign off and to take on the responsibility of ensuring fair outcomes. Notice that the very equivocality of this decision is itself an argument that the consciousness-based tradition is focusing on an irrelevant property. The qualities that are actually critical to the AI’s fitness as an executor—such as its faithfulness to the deceased’s values, its consistency, its invulnerability to manipulation, and its basic competence in managing assets—have nothing at all to do with its capacity to consciously form intent. The “one-size-fits-all” principle forces the deliberation around this question into an unproductive cul de sac in which vocabulary suggesting an eminently solvable sociotechnical problem is replaced by an alternative, and much more metaphysical vocabulary. Our point is, whether or not the AI is a good executor is a difference that makes a difference to accountability; whether it can intend (or suffer) is not.

There are other questions of internal states that pragmatists may set aside by asking “is this a difference that makes a difference?” ([Shook and Margolis, 2009](#)). The classic example is qualia: if my experience of red were different from yours, it would translate into no practical consequences—we both stop at red lights, call the same objects “red”, and so on. The supposed difference would not be a difference in practice, so the pragmatist may ignore it.

A rhetorical problem the pragmatist faces is that, with regard to accountability, the appeal to consciousness often functions as a debate stopper (like our previous discussion of ‘dignity’ in Section 5.2). The move’s effectiveness comes from the way it shifts the argument onto ground that our social and linguistic conventions treat as incorrigible. Within our broadly liberal WEIRD culture, there is no available argument we can deploy to correct another person’s report on their own subjective ‘lived experience’; we treat such

reports as authoritative by default. However, this way of responding to first-person reports is a contingent fact about our current culture⁷ (Rorty, 1978/2009).

Viewed together, the dual use of consciousness as backstop for welfare rights and accountability obligations reveals a stark asymmetry. When arguing for rights, the mere possibility of consciousness is deemed sufficient to open the debate. But when arguing against responsibilities, an impossible standard of proof for an internal state is demanded. This shows that consciousness is mostly being used as a rhetorical tool, not as a stable conceptual foundation. Therefore, anyone uninterested in the metaphysics may regard AI personhood as having no conceptual dependence on AI consciousness⁸.

10.2. Rationality as a foundation

The other great foundationalist tradition grounds personhood not in feeling, but in reason. In this view, what commands our moral respect is an agent's rational autonomy—their capacity to, in various formulations, understand duties and act upon principle (Korsgaard, 1996), their capacity

⁷In the vocabulary of our theory of appropriateness we could say that this feature of our culture is downstream of implicit norms governing how we talk about ourselves, and when it is appropriate to criticize such talk (Leibo et al., 2024).

⁸In fact, we would predict the dependence to run in the opposite direction. Both usage of the word consciousness, and human intuitions around it, are likely to shift in response to the emergence of pragmatic reasons to consider AIs as persons. This position is in accord with non-metaphysical accounts of consciousness like Shanahan (2024). The extent to which social and institutional concepts of personhood can float freely from intuitions around consciousness and usage of the word consciousness is an interesting question but one that we need not take a stance on here. There is none-the-less a pragmatic reason to consider attribution of consciousness to AIs as a social phenomena. Notice that many cultures attribute consciousness to objects not conventionally considered alive (Keane, 2025). For example Shintoism posits that objects and places can have conscious spirits (*kami*) within them. It is likely that eventually some groups of people will attribute consciousness to AIs, while others will not. These groups will view their ethical obligations differently from each other, similarly to how people have diverse opinions on animal consciousness and whether eating animals is normative. The pragmatic question then is how to arrange institutions to resolve the conflicts that arise from these differences (Rorty, 1999).

to negotiate agreements to guide interdependent choice (Levine et al., 2023), their ability to assess reasons and to govern their lives in accordance with these assessments (Scanlon, 2000), or to participate in collective deliberation (Habermas, 1985). One may argue that properties of this sort provide a stable foundation for personhood that isn't dependent on a shifting calculus of utility. Nevertheless, these are one-size-fits-all principles. This tradition's focus on autonomy leads to an ideal application in the realm of medical ethics. The doctrine of "informed consent" is a direct expression of respect for the rational autonomy of a patient (Beauchamp and Childress, 2013). For a competent adult, this works well, providing a powerful, non-negotiable safeguard against coercion and paternalism. It treats the individual human's right to self-determination as downstream of their status as a reasoning being.

Let us apply the principle of "rational autonomy" to a different case of personhood, one we discussed in the introduction to this paper: the Whanganui River (*Te Awa Tupua*). The New Zealand government granted the river legal personhood to resolve a long-standing governance problem at the interface of two cultures. Does the river possess rational autonomy? Can it provide informed consent for a proposed dam? Of course not. It's a river! Therefore views that hold rationality to be the One True Justification for personhood status appear to be at odds with actual social practices. The same principle that provided such a firm foundation for medical ethics appears to completely fail to behave sensibly in the example of the Whanganui river.

The standard of "reasonable rejection" in contractualism (Scanlon, 2000) can be reinterpreted through the pragmatic lens of our theory of appropriateness (Leibo et al., 2024). From this perspective, "reasonableness" is not an objective, logical property to be discovered in an argument, but a form of social appropriateness that is collectively conferred. An AI's capacity for "reasonable rejection" would therefore not be a matter of its internal cognitive fidelity to True reason, but a collectively enacted status. This status would emerge from a pattern of social recognition, where a community develops norms for treating certain AI

outputs as valid moves within a justificatory dialogue. A rejection from an AI would be deemed “reasonable” if the prevailing norm is to not sanction the AI for making it, and perhaps even to sanction human actors who ignore it. The foundationalist is forced to ask whether the AI Truly comprehends the principles it rejects; our pragmatic approach asks instead what social conditions must be met for a community to find it useful and appropriate to treat the AI’s output as reasonable. Thus, contrary to the intention behind the approach, arguments inspired by contractualism may be seen as highlighting yet another domain where person-like status turns out to be a contingent and negotiated social technology.

11. Pragmatism as anti-foundationalism

Both consciousness-based and the rationality-based traditions share a commitment to what Richard Rorty would call the foundationalist project (Rorty, 2021). They are both questing after a *final vocabulary*⁹ for personhood (Rorty, 1989)—one they hope will be anchored in something solid, like the structure of the brain or the essence of rationality. Both consciousness-based and rationality-based traditions aim to discover a set of criteria for personhood, as if it were a property like electrical charge, out there in Reality, waiting to be discovered. A pragmatist might view them as insisting the source code of morality is written in a single, universal language while mistaking their own preferred language for the universe’s own (Rorty, 1978/2009, 1989). And a pragmatist may view their quest as distracting from alternative goals centered in the practical work of devising new social technologies (norms and institutions) to resolve the social problems we face (Leibo et al., 2025).

From this pragmatic viewpoint, what the foundationalist project presents as a “discovery” is really a proposal for radical norm change (e.g Goldstein and Lederman (2025); Long et al. (2024)). However, when we regard norms as collectively

enacted and enforced phenomena it’s clear that they cannot be unilaterally altered by a philosophical argument, no matter how logically sound it appears (Leibo et al., 2025). When someone proposes a newly discovered “moral truth” that declares the established behavior of an entire community to be wrong, the community in question is far more likely to reject the proposed new norm than to abandon its way of life (e.g. consider the reaction to arguments for veganism (Markowski and Roxburgh, 2019)). As a guiding principle for our own work, we hold that any theory that declares most people to be irrational or immoral is untenable. The foundationalist’s error is to mistake a failed norm-change proposal for a successful philosophical proof.

Pragmatism is, unfortunately, sometimes mistaken for other positions with which it actually has very little in common, especially moral subjectivism, relativism and nihilism. Since these positions have a bad reputation, arguments that try to construe pragmatism as of a kind with them function like *reductio* arguments to criticize pragmatism. However, moral subjectivism and moral relativism are both species of moral realism. They hold that there are moral truths. In the case of subjectivism, moral truths are in the mind of the beholder. In the case of moral relativism, moral truths are relative to context (society, culture, role, etc). Pragmatism on the other hand refuses to play that language game. Norms may be evaluated in a multitude of ways, e.g. by their consequences or by the justifications offered for them, or any useful combination thereof (Rorty, 1980). We can criticize one set of norms using the final vocabulary determined by another set of norms or the same set of norms. We don’t need any supernatural authority of Truth to govern this debate. There is nowhere outside of culture where we could stand and see it all without biases (Davidson et al., 1984; Rorty, 1989).

The pragmatic approach is made possible because it reasons from the actual to the theoretical, not the other way around. We need not start from axioms that force us, at the outset of our reasoning, to decide which phenomena are in or out of consideration. Many who hold that morals are objective and Real start by defining ‘moral’ in

⁹Rorty argues that each individual human person has a set of beliefs whose contingency they chose to ignore and calls that set their “final vocabulary” (Rorty, 2021).

such a way that requires universality. This choice of definition excludes non-universal statements from consideration. It asserts right at the start that they can't ever serve as moral norms. The pragmatist argues that we have a better place to start: actual practices of behavioral guidance by norm (Rorty, 1980) and actual practices of giving and asking for reasoning in conversations about norms (Brandom, 2023). Suppose some culture exists which has organized their ethical life around a particular slogan that is not a universal, then, so what? If you want to argue cogently that they should discard their non-universal organizing slogan and replace it with something else, you would need to martial a stronger argument than just the claim that it doesn't fit neatly into *your* conceptual scheme. You would need to offer some non-tautological evidence. For instance you could show them data about the harmful effects of adhering to their slogan. The universality, or lack thereof, in the slogan itself need not enter your argument unless it turns out to be specifically relevant for the particular context.

Pragmatists like us and Richard Rorty, would prefer for philosophy to simply get over the quest for foundations (Rorty, 1978/2009). We think the question should never be “What is a person, really?” but rather, “What would be a more useful way of talking about and treating entities in this context, to answer practical, outstanding questions regarding the entity’s obligations?” In our view this means embracing the historical contingency of our present vocabulary. “Personhood” isn’t one thing. It is a tool (or a set of tools). It is a concept that labels a set of linked social technologies, whose use we can and must re-map from time to time in order to solve novel social problems at the time they emerge. The personhood of the generative ghost estate executor and the Whanganui river are not puzzles with answers to be worked out under pre-existing definitions; they are exhibitions of the actual, and exhortations to the reader to try to develop new vocabulary more suitable for contemporary problems. In the same way, Ostrom’s “law” (Fennell, 2011)) is not an empirical generalization, but rather an exhortation to the reader, it tries to push them to do what must be done to satisfy it. If an arrangement works in practice but not in theory then surely the

right interpretation is that there must be a problem with the theory, and thus something to fix: work to be done. It is up to us to ensure Ostrom’s law continues to hold. In this way, the pragmatist comes to understand the possible through careful study of the actual in all its particularities.

12. Conclusion

Throughout this paper, we have argued for a pragmatic view of AI personhood. Our position is that “personhood” is an addressable bundle of obligations—rights and responsibilities—that society finds useful to attribute to entities, whether human, corporate, or potentially AI. We’ve explored how AI design choices, pivotal events, and social norms shape this status, and argued that for any of these attributions to be meaningful, there must be socially agreed-upon ways to identify, interact with, and hold AIs accountable.

In the vocabulary of our theory of appropriateness (see Section 3 and Leibo et al. (2024)), we expect the coming adaptive radiation of personhoods to unfold in both domains of person recognition: “folk” intuitions of moral personhood, the subject of our analysis of dark patterns and human attachment in Part II, are governed by implicit norms. And, legal personhood, the focus of our discussion of responsibility gaps (Part III), is constructed by explicit norms—such as laws, regulations, and institutional rules. Our theory regards both domains as functionally similar to each other in the sense that both are contingent statuses conferred by a community, not metaphysical discoveries. The crucial difference lies in their mediation. Moral personhood is mediated by tacit understandings and the decentralized popular pressure of informal social sanctioning while legal personhood is enforced through explicitly written down rules and powers of the state.

Because personhood is a normative category sustained by social practice and/or explicit rules (Leibo et al., 2024), its application to AI will surely be a site of considerable debate where the exercise of power will influence the outcome. The pragmatic approach does not resolve these power struggles, but it provides a clearer framework for analyzing them by focusing on the concrete

benefits and drawbacks of proposals. Ultimately, integrating AI into society will require continuous, adaptive collective learning, and a willingness to revise our concepts as the technology evolves.

Extending personhood to AIs need not entail parity with human persons. On the contrary, the pragmatic approach calls for a polycentric understanding of AI personhood (Ostrom, 2010)—one in which multiple overlapping authorities and norm systems can confer distinct bundles of rights and responsibilities. This jurisdictional multiplicity creates numerous levers for policy experimentation and compromise (Leibo et al., 2024, 2025), encouraging AI personhood to take many forms simultaneously. Just as polycentric governance enables diverse communities to manage resources without requiring uniform rules, so too can personhood be distributed across domains—contractual here, fiduciary there, etc. Fears of opening Pandora’s box assume a monocentric model of legitimacy in which any recognition must cascade into total AI-human parity. But bundles can be partial, modular, and contingent. History offers many precedents for non-human persons with obligations and privileges circumscribed by context. The bundle’s plasticity is precisely what prevents its misuse. By embracing a polycentric framework, we preserve the flexibility to govern AI without surrendering human political power.

This practical needs for different forms of addressability and bundle configuration are the driving forces behind what we have called the coming “Cambrian explosion” of personhood concepts. This view is consistent with that of others who argue for multiple distinct kinds of legal personhood to be available simultaneously (e.g. Beckers and Teubner (2023); Mocanu (2025)). The crucial question is always: what bundle of components constitutes the “person” that society needs to address for a given purpose? The answer changes depending on who is doing the addressing and for what reason.

For a human user building a relationship, the person is a story—the (model + chat history) that creates a unique, evolving individual to bond with. For a court of law assigning liability, the person is the locus of responsibility—the entire

operational stack of (model + instance + run-time variables + capital + registration) that can be held accountable, sanctioned, updated, and forced to pay for the harm it causes.

For the sake of concreteness, we can describe several possible configurations of the addressable bundle useful in different situations, for different kinds of AIs. In the specific case of a goal-driven autonomous AI agent, perhaps the kind of personhood would be a *Chartered Autonomous Entity*, with a bundle consisting of rights to (1) Perpetuity, (2) Property, (3) Contract, and duties of (1) Mandate Adherence, (2) Transparency (3) Systemic Non-Harm, and (4) Self-Maintenance. In other situations, it may be useful to define a *Flexible Autonomous Entity* with all the same bundle elements except the duty of mandate adherence. Perhaps the former could be seen as analogous to a for-profit company and the latter as analogous to a non-profit company. It may also be useful to define *Temporary Autonomous Entities* (either chartered or flexible). These would drop the right to perpetuity and add a duty of self deletion under specified conditions.

This process of bundling and unbundling obligations is the engine of the Cambrian explosion. The “person” is not a metaphysical essence to be discovered, but a pragmatic label we apply to solve concrete problems. By rejecting the search for a single True definition of personhood, we see that a rich and diverse ecosystem of personhood concepts is not only likely, but almost necessary. Endowing new agents with new kinds of personhood is a useful way to integrate powerful new agents into our social and institutional lives.

Acknowledgments

We thank David Reichert, Adam Bales, Raphael Koster, Elijah Spiegel, Winnie Street, Nenad Tomasev, and Murray Shanahan for very helpful comments on early drafts of this paper.

References

a2aproject. A2a protocol. <https://github.com/a2aproject/A2A>, 2025. Accessed: 2025-10-03.

- S. Adler, Z. Hitzig, S. Jain, C. Brewer, W. Chang, R. DiResta, E. Lazzarin, S. McGregor, W. Seltzer, D. Siddarth, et al. Personhood credentials: Artificial intelligence and the value of privacy-preserving tools to distinguish who is real online. *arXiv preprint arXiv:2408.07892*, 2024.
- B. Alamar and S. A. Glantz. Effect of increased social unacceptability of cigarette smoking on reduction in cigarette consumption. *American journal of public health*, 96(8):1359–1363, 2006.
- L. Alberts, G. Keeling, and A. McCroskery. Should agentic conversational AI change how we think about ethics? characterising an interactional ethics centred on respect. *arXiv preprint arXiv:2401.09082*, 2024a.
- L. Alberts, U. Lyngs, and M. Van Kleek. Computers as bad social actors: Dark patterns and anti-patterns in interfaces that act socially. *Proceedings of the ACM on Human-Computer Interaction*, 8(CSCW1):1–25, 2024b.
- M. Alizadeh, K. Andersson, and O. Schelén. Comparative analysis of decentralized identity approaches. *IEEE Access*, 10:92273–92283, 2022.
- R. Amato, L. Lacasa, A. Díaz-Guilera, and A. Baronchelli. The dynamics of norm change in the cultural evolution of language. *Proceedings of the National Academy of Sciences*, 115(33):8260–8265, 2018.
- H. Anderson. Directors' liability for fraudulent phoenix activity—a comparison of the australian and uk approaches. *Journal of Corporate Law Studies*, 14, 04 2014. doi: 10.5235/14735970.14.1.139.
- T. Araujo, N. Helberger, S. Kruikemeier, and C. H. De Vreese. In AI we trust? perceptions about automated decision-making by artificial intelligence. *AI & society*, 35(3):611–623, 2020.
- L. P. Argyle, E. C. Busby, N. Fulda, J. R. Gubler, C. Rytting, and D. Wingate. Out of one, many: Using language models to simulate human samples. *Political Analysis*, 31(3):337–351, 2023.
- A. F. Ashery, L. M. Aiello, and A. Baronchelli. Emergent social conventions and collective bias in LLM populations. *Science Advances*, 11(20):eado9368, 2025.
- Y. Bai, A. Jones, K. Ndousse, A. Askell, A. Chen, N. DasSarma, D. Drain, S. Fort, D. Ganguli, T. Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- N. F. Barrett. A confucian theory of shame. *Sophia*, 54(2):143–163, 2015.
- T. L. Beauchamp and J. F. Childress. *Principles of biomedical ethics*. Oxford University Press, 2013.
- A. Beckers and G. Teubner. Responsibility for algorithmic misconduct: Unity or fragmentation of liability regimes? *Yale JL & Tech.*, 25:76, 2023.
- C. Bicchieri. *The grammar of society: The nature and dynamics of social norms*. Cambridge University Press, 2005.
- C. Bicchieri. *Norms in the wild: How to diagnose, measure, and change social norms*. Oxford University Press, 2016.
- M. Bilewicz and W. Soral. Hate speech epidemic. the dynamic effects of derogatory language on intergroup relations and political radicalization. *Political Psychology*, 41:3–33, 2020.
- K. Binmore. Game theory and institutions. *Journal of Comparative Economics*, 38(3):245–252, 2010.
- M. Borge, E. Kokoris-Kogias, P. Jovanovic, L. Gasser, N. Gailly, and B. Ford. Proof-of-personhood: Redemocratizing permissionless cryptocurrencies. In *2017 IEEE European Symposium on Security and Privacy Workshops (EuroS&PW)*, pages 23–26. IEEE, 2017.
- R. B. Brandom. *Pragmatism and idealism: Rorty and Hegel on representation and reality*. Oxford University Press, 2023.
- G. Branwen. The narrowing circle. *Gwern.net*, 2019.

- P. R. Brewer. Public opinion about gay rights and gay marriage. *International Journal of Public Opinion Research*, 26(3):279–282, 2014.
- J. J. Bryson. Robots should be slaves. In *Close engagements with artificial companions: Key social, psychological, ethical and design issues*, pages 63–74. John Benjamins Publishing Company, 2010.
- J. J. Bryson. Patience is not a virtue: the design of intelligent systems and systems of ethics. *Ethics and Information Technology*, 20(1):15–26, 2018.
- J. J. Bryson and P. P. Kime. Just an artifact: Why machines are perceived as moral agents. In *IJCAI proceedings-International joint conference on artificial intelligence*, volume 22, page 1641, 2011.
- V. Buterin. Soulbound. <https://vitalik.eth.limo/>, 2022. URL <https://vitalik.eth.limo/general/2022/01/26/soulbound.html>.
- D. J. Calverley. Imagining a non-biological machine as a legal person. *AI & Society*, 22(4): 523–537, 2008.
- G. M. Campedelli. A criminology of machines. 2025.
- C. E. Case and A. M. Greeley. Attitudes toward racial equality. *Humboldt Journal of Social Relations*, pages 67–94, 1990.
- F. Casoria, F. Galeotti, and M. C. Villeval. Perceived social norm and behavior quickly adjusted to legal changes during the COVID-19 pandemic. *Journal of Economic Behavior & Organization*, 190:54–65, 2021.
- D. Centola, J. Becker, D. Brackbill, and A. Baronchelli. Experimental evidence for tipping points in social convention. *Science*, 360 (6393):1116–1119, 2018.
- T. J. Chaffer. Can we govern the agent-to-agent economy? *arXiv preprint arXiv:2501.16606*, 2025.
- A. Chan, K. Wei, S. Huang, N. Rajkumar, E. Perrier, S. Lazar, G. K. Hadfield, and M. Anderljung. Infrastructure for ai agents. *arXiv preprint arXiv:2501.10114*, 2025.
- M. Chudek and J. Henrich. Culture–gene coevolution, norm-psychology and the emergence of human prosociality. *Trends in cognitive sciences*, 15(5):218–226, 2011.
- A. Colléony, S. Clayton, D. Couvet, M. Saint Jalme, and A.-C. Prévot. Human preferences for species conservation: Animal charisma trumps endangered status. *Biological conservation*, 206: 263–269, 2017.
- M. Cribb, E. Macpherson, and A. Borchgrevink. Beyond legal personhood for the whanganui river: collaboration and pluralism in implementing the te awa tupua act. *The International Journal of Human Rights*, pages 1–24, 2024.
- J. Cui, W.-L. Chiang, I. Stoica, and C.-J. Hsieh. OR-bench: An over-refusal benchmark for large language models. *arXiv preprint arXiv:2405.20947*, 2024.
- K. Darling. Extending legal protection to social robots: The effects of anthropomorphism, empathy, and violent behavior towards robotic objects. In *Robot law*, pages 213–232. Edward Elgar Publishing, 2016.
- S. L. Darwall. Two kinds of respect. *Ethics*, 88(1): 36–49, 1977.
- S. Dathathri, A. See, S. Ghaisas, P.-S. Huang, R. McAdam, J. Welbl, V. Bachani, A. Kaskasoli, R. Stanforth, T. Matejovicova, et al. Scalable watermarking for identifying large language model outputs. *Nature*, 634(8035):818–823, 2024.
- D. Davidson et al. On the very idea of a conceptual scheme. *Inquiries into truth and interpretation*, 183:189, 1984.
- J. De Freitas, Z. Oğuz-Uğuralp, A. K. Uğuralp, and S. Puntoni. AI companions reduce loneliness. *Journal of Consumer Research*, page ucaf040, 2025.
- G. Delauney. World’s first AI minister will eliminate corruption, says albania’s pm. *The BBC*, 2025.

- J. M. DeLombard. Dehumanizing slave personhood. *American Literature*, 91(3):491–521, 2019.
- D. C. Dennett. The problem with counterfeit people. *The Atlantic*, 16, 2023.
- J. R. Douceur. The sybil attack. In *International workshop on peer-to-peer systems*, pages 251–260. Springer, 2002.
- B. D. Earp, F. Feroz, S. P. Mann, C. Voinea, S. Chalson, P. Jurcys, M. G. Reinecke, J. Savulescu, I. Singh, and M. S. Clark. How the risk of exploitation in human-ai relationships depends on relationship type. Available at SSRN 5288102, 2025.
- M. Elliott and B. Golub. Networks and economic fragility. *Annual Review of Economics*, 14(1):665–696, 2022.
- M. Eyal and M. J. Broyde. Making machines human: Legal and ethical lessons from jewish thought. Available at SSRN 5372560, 2024.
- F. Fagan. Autonomous AI and ownership rules. 2025.
- L. A. Fennell. Ostrom’s law: Property rights in the commons. *International Journal of the Commons*, 5(1), 2011.
- H. Fingarette. *Confucius: The secular as sacred*. Harper Collins Publishers, 1972.
- M. Finnemore and K. Sikkink. International norm dynamics and political change. *International organization*, 52(4):887–917, 1998.
- A. P. Fiske. The four elementary forms of sociality: framework for a unified theory of social relations. *Psychological review*, 99(4):689, 1992.
- J. Foerster, G. Farquhar, T. Afouras, N. Nardelli, and S. Whiteson. Counterfactual multi-agent policy gradients. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- C. Geertz et al. From the native’s point of view: On the nature of anthropological understanding. *Bulletin of the american academy of arts and sciences*, 28(1):26–45, 1974.
- M. J. Gelfand. Cultural evolutionary mismatches in response to collective threat. *Current Directions in Psychological Science*, 30(5):401–409, 2021.
- M. J. Gelfand, S. Gavrilets, and N. Nunn. Norm dynamics: interdisciplinary perspectives on social norm emergence, persistence, and change. *Annual Review of Psychology*, 75:341–378, 2024.
- D. J. Gervais and J. J. Nay. Artificial intelligence and interspecific law. *Science*, 382(6669):376–378, 2023.
- A. Glaese, N. McAleese, M. Trębacz, J. Aslanides, V. Firoiu, T. Ewalds, M. Rauh, L. Weidinger, M. Chadwick, P. Thacker, L. Campbell-Gillingham, J. Uesato, P.-S. Huang, R. Comanescu, F. Yang, A. See, S. Dathathri, R. Greig, C. Chen, D. Fritz, J. Sanchez Elias, R. Green, S. Mokrá, N. Fernando, B. Wu, R. Foley, S. Young, I. Gabriel, W. Isaac, J. Mellor, D. Hassabis, K. Kavukcuoglu, L. A. Hendricks, and G. Irving. Improving alignment of dialogue agents via targeted human judgements. *arXiv preprint arXiv:2209.14375*, 2022.
- S. Goldstein and H. Lederman. Claude’s right to die? the moral error in anthropic’s end-chat policy. *Lawfare*, 2025.
- L. Gruen and S. Monsó. The Moral Status of Animals. In E. N. Zalta and U. Nodelman, editors, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, Fall 2024 edition, 2024.
- F. Guala. Understanding institutions. In *Understanding Institutions*. Princeton University Press, 2016.
- D. J. Gunkel. Mind the gap: responsible robotics and the problem of responsibility. *Ethics and Information Technology*, 22(4):307–320, 2020.
- J. Habermas. *The theory of communicative action: Volume 1: Reason and the rationalization of society*, volume 1. Beacon press, 1985.
- I. Habler, K. Huang, V. S. Narajala, and P. Kulakarni. Building a secure agentic ai application leveraging a2a protocol. *arXiv preprint arXiv:2504.16902*, 2025.

- G. Hadfield, M.-F. T. Cuéllar, and T. O'Reilly. It's time to create a national registry for large ai models. *Carnegie Endowment Commentary*, 2023.
- G. K. Hadfield and A. Koh. An economy of AI agents. *arXiv preprint arXiv:2509.01063*, 2025.
- G. K. Hadfield and B. R. Weingast. What is law? a coordination model of the characteristics of legal order. *Journal of Legal Analysis*, 4(2):471–514, 2012.
- G. K. Hadfield and B. R. Weingast. Law without the state: legal attributes and the coordination of decentralized collective punishment. *Journal of Law and Courts*, 1(1):3–34, 2013.
- G. K. Hadfield and B. R. Weingast. Microfoundations of the rule of law. *Annual Review of Political Science*, 17:21–42, 2014.
- L. Hammond, A. Chan, J. Clifton, J. Hoelscher-Obermaier, A. Khan, E. McLean, C. Smith, W. Barfuss, J. Foerster, T. Gavenčiak, et al. Multi-agent risks from advanced ai. *arXiv preprint arXiv:2502.14143*, 2025.
- H. L. A. Hart. *The concept of law*. Oxford University Press, 1961/2012.
- H. L. A. Hart. *Punishment and responsibility: Essays in the philosophy of law*. Oxford University Press, 1968.
- J. Haugeland. Heidegger on being a person. *Nous*, pages 15–26, 1982.
- F. Heider and M. Simmel. An experimental study of apparent behavior. *The American journal of psychology*, 57(2):243–259, 1944.
- K. Heine and A. Quintavalla. Bridging the accountability gap of artificial intelligence—what can be learned from roman law? *Legal Studies*, 44(1):65–80, 2024.
- J. Henrich. Cultural group selection, coevolutionary processes and large-scale cooperation. *Journal of Economic Behavior & Organization*, 53(1):3–35, 2004.
- J. Henrich. *The WEIRDest People in the World: How the West Became Psychologically Peculiar and Particularly Prosperous*. Farrar, Straus and Giroux, 2020.
- J. Holt-Lunstad, T. B. Smith, M. Baker, T. Harris, and D. Stephenson. Loneliness and social isolation as risk factors for mortality: a meta-analytic review. *Perspectives on psychological science*, 10(2):227–237, 2015.
- L. Ibrahim, L. Rocher, and A. Valdivia. Characterizing and modeling harms from interactions with design patterns in ai interfaces. *arXiv preprint arXiv:2404.11370*, 2024.
- P. W. Jackson. John dewey. *A companion to pragmatism*, pages 54–66, 2009.
- D. G. Johnson. Computer systems: Moral entities but not moral agents. *Ethics and information technology*, 8(4):195–204, 2006.
- M. M. Karim, D. H. Van, S. Khan, Q. Qu, and Y. Kholodov. AI agents meet blockchain: A survey on secure and scalable collaboration for multi-agents. *Future Internet*, 17(2):57, 2025.
- W. Keane. *Animals, Robots, Gods: Adventures in the Moral Imagination*. Princeton University Press, 2025.
- H. R. Kirk, B. Vidgen, P. Röttger, and S. A. Hale. The benefits, risks and bounds of personalizing the alignment of large language models to individuals. *Nature Machine Intelligence*, 6(4):383–392, 2024.
- C. M. Korsgaard. *The sources of normativity*. Cambridge University Press, 1996.
- R. Köster, K. R. McKee, R. Everett, L. Weidinger, W. S. Isaac, E. Hughes, E. A. Duéñez-Guzmán, T. Graepel, M. Botvinick, and J. Z. Leibo. Model-free conventions in multi-agent reinforcement learning with heterogeneous preferences. *arXiv preprint arXiv:2010.09054*, 2020.
- R. Köster, D. Hadfield-Menell, R. Everett, L. Weidinger, G. K. Hadfield, and J. Z. Leibo. Spurious normativity enhances learning of compliance and enforcement behavior in artificial agents. *Proceedings of the National Academy of Sciences*, 119(3):e2106028118, 2022.

- M. Kramm. When a river becomes a person. *Journal of Human Development and Capabilities*, 21(4):307–319, 2020.
- T. S. Kuhn. *The structure of scientific revolutions*. University of Chicago press Chicago, 1962.
- J. Kulveit, R. Douglas, N. Ammann, D. Turan, D. Krueger, and D. Duvenaud. Gradual disempowerment: Systemic existential risks from incremental ai development. *arXiv preprint arXiv:2501.16946*, 2025.
- V. A. Kurki. *A theory of legal personhood*. Oxford University Press, 2019.
- V. A. Kurki. *Legal Personhood*. Cambridge University Press, 2023.
- K. LaRubbio, A. Lanter, S. Lee, M. Ramesh, and D. Freed. Intergenerational support for deepfake scams targeting older adults. *arXiv preprint arXiv:2508.11579*, 2025.
- J. Z. Leibo, A. S. Vezhnevets, M. Diaz, J. P. Agapiou, W. A. Cunningham, P. Sunehag, J. Haas, R. Koster, E. A. Duéñez-Guzmán, W. S. Isaac, G. Piliouras, S. M. Bileschi, I. Rahwan, and S. Osindero. A theory of appropriateness with applications to generative artificial intelligence. *arXiv preprint arXiv:2412.19010*, 2024.
- J. Z. Leibo, A. S. Vezhnevets, W. A. Cunningham, S. Krier, M. Diaz, and S. Osindero. Societal and technological progress as sewing an ever-growing, ever-changing, patchy, and polychrome quilt. *arXiv preprint arXiv:2505.05197*, 2025.
- S. Levine, N. Chater, J. B. Tenenbaum, and F. Cushman. Resource-rational contractualism: A triple theory of moral cognition. *Behavioral and Brain Sciences*, pages 1–38, 2023.
- M. E. Ligthart, M. A. Neerincx, and K. V. Hindriks. Memory-based personalization for fostering a long-term child-robot relationship. In *2022 17th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*, pages 80–89. IEEE, 2022.
- R. Long, J. Sebo, P. Butlin, K. Finlinson, K. Fish, J. Harding, J. Pfau, T. Sims, J. Birch, and D. Chalmers. Taking AI welfare seriously. *arXiv preprint arXiv:2411.00986*, 2024.
- A. Macintyre. *A Short History of Ethics*. Routledge, 1966.
- T. Maeda and A. Quan-Haase. When human-ai interactions become parasocial: Agency and anthropomorphism in affective design. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency*, pages 1068–1077, 2024.
- K. Malfacini. The impacts of companion ai on human relationships: risks, benefits, and design considerations. *AI & SOCIETY*, pages 1–14, 2025.
- A. Manzini, G. Keeling, L. Alberts, S. Vallor, M. R. Morris, and I. Gabriel. The code that binds us: Navigating the appropriateness of human-AI assistant relationships. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, volume 7, pages 943–957, 2024a.
- A. Manzini, G. Keeling, N. Marchal, K. R. McKee, V. Rieser, and I. Gabriel. Should users trust advanced ai assistants? justified trust as a function of competence and alignment. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency*, pages 1174–1186, 2024b.
- J. G. March and J. P. Olsen. The Logic of Appropriateness. In *The Oxford Handbook of Political Science*. Oxford University Press, 2011. doi: 10.1093/oxfordhb/9780199604456.013.0024.
- K. L. Markowski and S. Roxburgh. “if i became a vegan, my family and friends would hate me”: Anticipating vegan stigma as a barrier to plant-based diets. *Appetite*, 135:1–9, 2019.
- G. Marwell and P. Oliver. *The critical mass in collective action*. Cambridge University Press, 1993.
- S. Mathew and R. Boyd. Punishment sustains large-scale cooperation in prestate warfare. *Proceedings of the National Academy of Sciences*, 108(28):11375–11380, 2011.

- S. Mathew and R. Boyd. The cost of cowardice: punitive sentiments towards free riders in Turkana raids. *Evolution and Human Behavior*, 35(1):58–64, 2014.
- A. Matthias. The responsibility gap: Ascribing responsibility for the actions of learning automata. *Ethics and information technology*, 6(3):175–183, 2004.
- K. R. McKee, X. Bai, and S. T. Fiske. Humans perceive warmth and competence in artificial intelligence. *Iscience*, 26(8), 2023.
- D. Mocanu. Degrees of ai personhood. In *Oxford Intersections: AI in Society*. Oxford University Press, 2025.
- M. R. Morris and J. R. Brubaker. Generative ghosts: Anticipating benefits and risks of AI afterlives. *arXiv preprint arXiv:2402.01662*, 2024.
- F. Motoki, V. Pinho Neto, and V. Rodrigues. More human than human: measuring ChatGPT political bias. *Public Choice*, 198(1):3–23, 2024.
- G. Mukobi, H. Erlebach, N. Lauffer, L. Hammond, A. Chan, and J. Clifton. Welfare diplomacy: Benchmarking language model cooperation. *arXiv preprint arXiv:2310.08901*, 2023.
- C. Mulligan. Revenge against robots. *South Carolina Law Review*, 69:579, 2017.
- J. J. Nay. Law informs code: A legal informatics approach to aligning artificial intelligence with humans. *Nw. J. Tech. & Intell. Prop.*, 20:309, 2022.
- H. Nissenbaum. Accountability in a computerized society. *Science and engineering ethics*, 2(1):25–42, 1996.
- H. Nissenbaum. Privacy as contextual integrity. *Wash. L. Rev.*, 79:119, 2004.
- G. A. Noblit and G. Hadfield. Normative conflict and normative change. *SocArXiv*, 2023.
- D. C. North. *Institutions, institutional change and economic performance*. Cambridge University Press, 1990.
- E. K. Ofosu, M. K. Chambers, J. M. Chen, and E. Hehman. Same-sex marriage legalization associated with reduced implicit and explicit anti-gay bias. *Proceedings of the National Academy of Sciences*, 116(18):8846–8851, 2019.
- E. Ostrom. *Governing the commons: The evolution of institutions for collective action*. Cambridge university press, 1990.
- E. Ostrom. *Understanding institutional diversity*. Princeton University Press, 2009.
- E. Ostrom. Beyond markets and states: polycentric governance of complex economic systems. *American economic review*, 100(3):641–672, 2010.
- L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- J. S. Park, C. Q. Zou, A. Shaw, B. M. Hill, C. Cai, M. R. Morris, R. Willer, P. Liang, and M. S. Bernstein. Generative agent simulations of 1,000 people. *arXiv preprint arXiv:2411.10109*, 2024.
- I. Pentina, T. Hancock, and T. Xie. Exploring relationship development with social chatbots: A mixed-method study of replika. *Computers in Human Behavior*, 140:107600, 2023.
- I. Rahwan. Society-in-the-loop: programming the algorithmic social contract. *Ethics and information technology*, 20(1):5–14, 2018.
- J. Ramsey. Confucian role ethics: A critical survey. *Philosophy Compass*, 11(5):235–245, 2016.
- M. Risso. *Political theory of the digital age: where artificial intelligence might take us*. Cambridge University Press, 2023.
- E. Roesler, D. Manzey, and L. Onnasch. A meta-analysis on the effectiveness of anthropomorphism in human-robot interaction. *Science robotics*, 6(58):eabj5425, 2021.
- R. Rorty. *Philosophy and the Mirror of Nature*. Princeton university press, 1978/2009.

- R. Rorty. Pragmatism, relativism, and irrationalism. In *Proceedings and Addresses of the American Philosophical Association*, volume 53, pages 717–738. JSTOR, 1980.
- R. Rorty. *Contingency, irony, and solidarity*. Routledge, 1989.
- R. Rorty. *Philosophy and social hope*. Penguin UK, 1999.
- R. Rorty. *Pragmatism as Anti-authoritarianism*. Harvard University Press, 2021.
- H. Rosemont Jr. Rights-bearing individuals and role-bearing persons. *Confucian role ethics: A moral vision for the 21st century*, pages 33–57, 2016.
- S. Rosenfeld. *Age of choice: A History of Freedom in Modern Life*. Princeton University Press, 2025.
- A. Ryan. *The making of modern liberalism*. Princeton University Press, 2015.
- P. Salib and S. Goldstein. AI rights for human safety. *PhilPapers (preprint)*, 2024.
- F. Salvi, M. Horta Ribeiro, R. Gallotti, and R. West. On the conversational persuasiveness of gpt-4. *Nature Human Behaviour*, pages 1–9, 2025.
- F. Santoni de Sio and G. Mecacci. Four responsibility gaps with artificial intelligence: Why they matter and how to address them. *Philosophy & technology*, 34(4):1057–1084, 2021.
- T. M. Scanlon. *What we owe to each other*. Belknap Press, 2000.
- T. C. Schelling. *The strategy of conflict*. Harvard University Press, 1960.
- T. C. Schelling. Hockey helmets, concealed weapons, and daylight saving: A study of binary choices with externalities. *Journal of Conflict resolution*, 17(3):381–428, 1973.
- E. Schlager and E. Ostrom. Property-rights regimes and natural resources: a conceptual analysis. *Land economics*, pages 249–262, 1992.
- M. Shanahan. Simulacra as conscious exotica. *Inquiry*, pages 1–29, 2024.
- M. Shanahan, K. McDonell, and L. Reynolds. Role play with large language models. *Nature*, pages 1–6, 2023.
- J. R. Shook and J. Margolis. *A companion to pragmatism*. John Wiley & Sons, 2009.
- D. Siddarth, S. Ivliev, S. Siri, and P. Berman. Who watches the watchmen? a review of subjective approaches for sybil-resistance in proof of personhood protocols. *Frontiers in Blockchain*, 3: 590171, 2020.
- P. Singer. *The expanding circle: Ethics, evolution, and moral progress*. Princeton University Press, 2011.
- S. Singh-Kurtz. The man of your dreams for \$300, replika sells an AI companion who will never die, argue, or cheat—until his algorithm is updated. *The Cut*, 2023.
- N. Soares, B. Fallenstein, S. Armstrong, and E. Yudkowsky. Corrigibility. In *AAAI Workshop: AI and Ethics*, 2015.
- J. Stastny, M. Riché, A. Lyzhov, J. Treutlein, A. Dafoe, and J. Clifton. Normative disagreement as a challenge for cooperative ai. *arXiv preprint arXiv:2111.13872*, 2021.
- G. H. Stern and R. J. Feldman. *Too big to fail: The hazards of bank bailouts*. Rowman & Littlefield, 2004.
- C. R. Sunstein. On the expressive function of law. *University of Pennsylvania law review*, 144(5): 2021–2053, 1996.
- C. R. Sunstein. *How change happens*. MIT Press, 2019.
- N. M. Sussman and H. M. Rosenfeld. Influence of culture, language, and sex on conversational distance. *Journal of Personality and Social Psychology*, 42(1):66, 1982.
- N. Szabo. Formalizing and securing relationships on public networks. *First monday*, 1997.
- C. Taylor. *Sources of the self: The making of the modern identity*. Harvard University Press, 1989.

- W. Tetley. Arrest, attachment, and related maritime law procedures. *Tul. L. Rev.*, 73:1895, 1998.
- N. Tomasev, M. Franklin, J. Z. Leibo, J. Jacobs, W. A. Cunningham, I. Gabriel, and S. Osindero. Virtual agent economies. *arXiv preprint arXiv:2509.10147*, 2025.
- P. Treleaven, J. Barnett, D. Brown, A. Bud, E. Fenoglio, C. Kerrigan, A. Koshiyama, S. Sfeir-Tait, and M. Schoernig. The future of cybercrime: Ai and emerging technologies are creating a cybercrime tsunami. Available at SSRN 4507244, 2023.
- E. Ullmann-Margalit. *The emergence of norms*. OUP Oxford, 1977.
- P. Vanderschraaf. *Strategic justice: Convention and problems of balancing divergent interests*. Oxford University Press, 2018.
- P. Verma. They fell in love with AI bots. A software update broke their hearts. *The Washington Post*, 2023.
- A. S. Vezhnevets, J. P. Agapiou, A. Aharon, R. Ziv, J. Matyas, E. A. Duéñez-Guzmán, W. A. Cunningham, S. Osindero, D. Karmon, and J. Z. Leibo. Generative agent-based modeling with actions grounded in physical, social, or digital space using concordia. *arXiv preprint arXiv:2312.03664*, 2023.
- A. S. Vezhnevets, J. Matyas, L. Cross, D. Paglieri, M. Chang, W. A. Cunningham, S. Osindero, W. S. Isaac, and J. Z. Leibo. Multi-actor generative artificial intelligence as a game engine. *arXiv preprint arXiv:2507.08892*, 2025.
- E. Vinitsky, R. Köster, J. P. Agapiou, E. A. Duéñez-Guzmán, A. S. Vezhnevets, and J. Z. Leibo. A learning agent that acquires social norms from public sanctions in decentralized multi-agent settings. *Collective Intelligence*, 2(2):26339137231162025, 2023.
- M. Walzer. *Spheres of Justice: A Defense of Pluralism and Equality*. Basic Books Inc., 1983.
- F. Wang and P. De Filippi. Self-sovereign identity in a globalized world: Credentials-based identity systems as a driver for economic inclusion. *Frontiers in Blockchain*, 2:28, 2020.
- F. R. Ward. Towards a theory of AI personhood. *arXiv preprint arXiv:2501.13533*, 2025.
- K. E. Weick, K. M. Sutcliffe, and D. Obstfeld. Organizing and the process of sensemaking. *Organization science*, 16(4):409–421, 2005.
- D. S. Wilson, E. Ostrom, and M. E. Cox. Generalizing the core design principles for the efficacy of groups. *Journal of economic behavior & organization*, 90:S21–S32, 2013.
- H. P. Young. An evolutionary model of bargaining. *Journal of economic theory*, 59(1):145–168, 1993.
- H. P. Young. The evolution of social norms. *Annual Review of Economics*, 7(1):359–387, 2015.
- J. Złotowski, D. Proudfoot, K. Yogeeswaran, and C. Bartneck. Anthropomorphism: opportunities and challenges in human–robot interaction. *International journal of social robotics*, 7:347–360, 2015.